

LEARNING LOOPY GRAPHICAL MODELS WITH LATENT VARIABLES: EFFICIENT METHODS AND GUARANTEES

BY ANIMASHREE ANANDKUMAR* AND RAGUPATHYRAJ VALLUVAN*

Univ. of California Irvine

The problem of structure estimation in latent graphical models is considered, where some nodes are latent or hidden. We characterize conditions for tractable graph estimation and develop efficient methods with provable guarantees. We consider models where the underlying Markov graph is locally tree-like and the model is in the regime of correlation decay. For the special case of the Ising model, the number of samples n required for structural consistency of our method scales as $n = \Omega(\theta_{\min}^{-\delta\eta(\eta+1)-2} \log p)$, where θ_{\min} is the minimum edge potential, δ is the depth (i.e., distance from a hidden node to the nearest observed nodes), and η is a parameter which depends on the bounds on node and edge potentials in the Ising model. Necessary conditions for structural consistency under any algorithm are derived and our method nearly matches the lower bound on sample requirements. Further, the proposed method is practical to implement and provides flexibility to control the number of latent variables and the cycle lengths in the output graph.

1. Introduction. It is widely recognized that the process of fitting the observed samples to a statistical model needs to incorporate latent or hidden factors, which are not directly observed. Learning latent variable models from observed samples involves mainly two tasks: discovering relationships between the observed and hidden variables, and estimating the strength of such relationships.

One of the simplest latent variable model is the so-called *latent class model* or *naïve Bayes model*, where the observed variables are conditionally independent given the state of the latent factor. An extension of these models are the class of *latent tree models* with many hidden variables forming a tree hierarchy. Latent tree models have been effective in modeling data in a variety of domains, such as the evolutionary process which gave rise to the present-day species in bio-informatics (popularly known as *phylogenetic tree models*) [23, 48], for financial and topic modeling [19], and for modeling contextual information for object recognition in computer vision [18]. There has been extensive work on learning latent tree models (e.g. [19, 25, 39]), where it is demonstrated that latent tree models can be learnt efficiently in high dimensions. In other words, the number of samples required for consistent learning is much smaller than the number of variables at hand. Moreover, inference on latent tree models is computationally tractable and can be carried out using simple algorithms such as *belief propagation*.

Despite all the above advantages, latent tree models may not be suitable in all scenarios and the assumption of a tree structure may be too restrictive. For instance, when latent trees are used to model topic-word relationships, the hypothesis is that the topics (which are hidden) and words are Markov on a tree. In other words, latent tree models posit that the words are generated from a single topic, while, in reality there are common words across topics. Loopy graphical models

*Supported by the setup funds at UCI and the AFOSR Award FA9550-10-1-0310.

AMS 2000 subject classifications: Primary 62H12; secondary 05C12

Keywords and phrases: Graphical model selection, Latent variables, Quartet methods

are able to capture such relationships and we consider learning such models, while retaining many advantages of the latent tree models.

Relaxing the tree assumption leads to non-trivial challenges: in general, learning these models is NP-hard [9, 32], even when there are no latent variables, and developing methods for learning such fully observed models is itself an area of active research (e.g. [4, 31, 45]). In this paper, we consider structure estimation in latent graphical models Markov on *locally tree-like* graphs, meaning that local neighborhoods in the graph do not contain cycles. Learning such graphs has many non-trivial challenges: are there parameter regimes where these models can be learnt consistently and efficiently? If so, are there practical learning algorithms? Are learning guarantees for loopy models comparable to those for latent trees? How does learning depend on various graph attributes such as node degrees, girth of the graph, and so on? We provide answers to these questions in this paper.

1.1. Our Approach and Contributions. We consider learning latent graphical models Markov on locally tree-like graphs in the regime of correlation decay. In this regime, there are no long-range correlations, and the local statistics converge to a tree limit. The implication of correlation decay is immediately clear: we can employ the available latent tree methods to learn “local” subgraphs consistently, as long as they do not contain any cycles. However, a non-trivial challenge remains: how does one merge these estimated local subgraphs (i.e., latent trees) to obtain an overall graph estimate? Specifically, merging involves matching latent nodes across different latent tree estimates, and it is not clear if this can be performed in an efficient manner.

We employ a different philosophy for building locally tree-like graphs with latent variables. We decouple the process of introducing cycles and latent variables in the output model. We initialize a loopy graph consisting of only the observed variables, and then iteratively add latent variables to local neighborhoods of the graph. We establish correctness of our method under a set of natural conditions.

We provide precise conditions for structural consistency of **LocalCLGrouping** under the probably approximately correct (PAC) model of learning [33] for general discrete models. We simplify these conditions for the Ising model, where each node is a binary random variable, to obtain better intuitions. We establish that for structural consistency, the number of samples is required to scale as $n = \Omega(\theta_{\min}^{-\delta\eta(\eta+1)-2} \log p)$, where p is the number of observed variables, θ_{\min} is the minimum edge potential, δ is the depth (i.e., graph distance from a hidden node to the nearest observed nodes), and η is a parameter which depends on the minimum and maximum node and edge potentials of the Ising model ($\eta = 1$ for homogeneous models). When there are no hidden variables ($\delta = 1$), the sample complexity is strengthened to $n = \Omega(\theta_{\min}^{-2} \log p)$, which matches with the best known sample complexity for learning fully-observed Ising models [4, 31].

We also establish necessary conditions for any (deterministic) algorithm to recover the graph structure, and establish that $n = \Omega(\Delta_{\min} \rho^{-1} \log p)$ samples are necessary for structural consistency, where Δ_{\min} is the minimum degree and ρ is the fraction of observed nodes. This is comparable to the requirement of the proposed method under uniform node sampling (i.e., selecting the observed nodes uniformly), given by $n = \Omega(\Delta_{\max}^2 \rho^{-2} (\log p)^3)$. Thus, our method is competitive with respect to the lower bound on learning.

Our proposed method has a number of attractive features for practical implementation: the method is amenable to parallelization which makes it efficient on large datasets. The method provides flexibility to control the length of cycles and the number of latent variables introduced in the output model. The method can incorporate penalty scores such as the Bayesian information criterion (BIC) [47] to tradeoff model complexity and fidelity. Moreover, by controlling the cycle lengths in the output model, we can obtain models with good inference accuracy under simple algo-

gorithms such as loopy belief propagation (LBP). Preliminary experiments on the newsgroup dataset suggests that the method can discover intuitive relationships efficiently, and also compares well with the popular latent Dirichlet allocation (LDA) [8] in terms of topic coherence and perplexity.

1.2. Related Work. The classical *latent cluster models* (LCM) consists of multivariate distributions with a single latent variable and the observed variables are conditionally independent under each state of the latent variable [36]. Hierarchical latent class (HLC) models [17, 52, 53] generalize these models by allowing multiple latent variables. However, the proposed learning algorithms are based on greedy local search in a high-dimensional space, which is computationally expensive. Moreover, the algorithms do not have theoretical guarantees. Similar shortcomings also hold for expectation-maximization (EM) based approaches [24, 34]. Learning latent trees has been studied extensively before, mainly in the context of phylogenetics. See [23, 48] for a thorough overview. Efficient algorithms with provable performance guarantees are available (e.g. [2, 19, 21, 25]). Our proposed method in this paper is inspired by [19].

Works on high-dimensional graphical model selection are more recent. The approaches can be mainly classified into two groups: non-convex local approaches [4, 11, 31, 41] and those based on convex optimization [15, 37, 45, 46]. There is a general agreement that the success of these methods is related to the presence of correlation decay in the model [4, 6]. This work makes the connection explicit: it relates the extent of correlation decay (i.e., the convergence rate to the tree limit) with the learning efficiency for latent models on large girth graphs. An analogous study of the effect of correlation decay for learning fully observed models is presented in [4].

This paper is the first work to provide provable guarantees for learning discrete latent models on loopy graphs in high dimensions (which can also be easily be extended to Gaussian models, see remarks following Theorem 2). The work in [16] considers learning latent Gaussian graphical models using a convex relaxation method. However, the method cannot be easily extended to discrete models. Moreover, the “incoherence” conditions required for the success of convex methods are hard to interpret and verify in general. In contrast, our conditions for success are transparent and based on the presence of correlation decay in the model. The work in [11] considers graphical model selection with hidden variables, but proposes learning Markov graph of marginal distribution (upon marginalizing the hidden variables) and then replacing the cliques in the estimated graphs with hidden variables. Sample complexity results are not provided, and the method performs poorly in high dimensions, since it aims to estimate dense graphs.

In [1], the problem of network tomography on locally tree-like graphs is considered, where the task is to estimate the graph using end-to-end path-based measurements (e.g. delay, link loss rate). It is established that a decaying fraction of participants is sufficient to learn the underlying graph. This paper has a different model, which is more challenging since here, we do not have a simple additive metric along the paths in the graph.

2. System Model.

2.1. Graphical Models. A *graphical model* is a family of multivariate distributions which are Markov in accordance to a particular undirected graph [35]. Each node in the graph $i \in W$ is associated to a random variable X_i taking value in a set \mathcal{X} . We consider discrete graphical models where \mathcal{X} is a finite set. The set of edges E captures the set of conditional independence relations among the random variables. We say that a set of random variables $\mathbf{X}_W := \{X_i, i \in W\}$ with probability mass function (pmf) P is Markov on the graph G if the local Markov property

$$(1) \quad P(x_i | x_{N(i)}) = P(x_i | x_{W \setminus i})$$

holds for all nodes $i \in W$, where $\mathcal{N}(i)$ are the neighbors of node i in graph G . More generally, we say that P satisfies the global Markov property, if for all disjoint sets $A, B \subset W$, we have

$$(2) \quad P(\mathbf{x}_A, \mathbf{x}_B | \mathbf{x}_S) = P(\mathbf{x}_A | \mathbf{x}_S) P(\mathbf{x}_B | \mathbf{x}_S).$$

where the set S is a *separator*¹ between A and B . The local and global Markov properties are equivalent under the *positivity* condition, given by $P(\mathbf{x}_W) > 0$, for all $\mathbf{x}_W \in \mathcal{X}^{|W|}$ [35], and we consider such distributions.

The Hammersley-Clifford theorem [10] states that under the positivity condition, a distribution P satisfies the Markov property according to a graph G iff. it factorizes according to the cliques of G , and we can write it in the exponential form as

$$(3) \quad P(\mathbf{x}) = \exp \left(\sum_{c \in \mathcal{C}} \theta_c(\mathbf{x}_c) - A(\boldsymbol{\theta}) \right),$$

where \mathcal{C} is the set of cliques of G and \mathbf{x}_c is the set of random variables on clique c . The quantity $A(\boldsymbol{\theta})$ is known as the *log-partition function* and serves to normalize the probability distribution. The functions θ_c are known as *potential* functions and correspond to the *canonical* parameters of the exponential family.

A special case is the Ising model, which is the pairwise model over binary variables $\{-1, +1\}^p$ with the probability mass function (pmf) given by

$$(4) \quad P(\mathbf{x}_W) = \exp \left(\sum_{e \in E} \theta_{i,j} x_i x_j + \sum_{i \in V} \phi_i x_i - A(\boldsymbol{\theta}) \right).$$

We specialize some of our results to the class of Ising models.

We consider latent graphical models in which a subset of nodes is latent or hidden. Let $H \subset W$ denote the hidden nodes and $V \subset W$ denote the observed nodes. Our goal is to discover the presence of hidden variables \mathbf{X}_H and learn the unknown graph structure $G(W)$, given n i.i.d. samples from observed variables \mathbf{X}_V . Let $p := |V|$ denote the number of observed nodes and $m := |W|$ denote the total number of nodes.

2.2. Tractable Graph Families: Girth-Constrained Graphs. In general, structure estimation of graphical models is NP-hard [9, 32]. We now characterize a tractable class of models for which we can provide guarantees on graph estimation.

We consider the family of graphs with a bound on the *girth*, which is the length of the shortest cycle in the graph. Let $\mathcal{G}_{\text{Girth}}(m; g)$ denote the ensemble of graphs with girth at least g . There are many graph constructions which lead to a bound on girth. For example, the bipartite Ramanujan graph [20, p. 107] and the random Cayley graphs [28] have bounds on the girth. Recently, efficient algorithms have been proposed to generate large girth graphs efficiently [5].

Although girth-constrained graphs are locally tree-like, in general, their global structure makes them hard instances for learning. Specifically, girth-constrained graphs have a large tree-width: it is known that a graph with average degree at least Δ_{avg} and girth at least g has a tree width as $\Omega \left(\frac{1}{g+1} (\Delta_{\text{avg}} - 1)^{\lfloor (g-1)/2 \rfloor} \right)$ [14]. Thus, learning is non-trivial for graphical models Markov on girth-constrained graphs, even when there are no latent variables due to their large treewidth [32].

¹A set $S \subset W$ is a separator for sets A and B if the removal of nodes in S separates A and B into distinct components.

2.3. Local Convergence to a Tree Limit. This work establishes tractable learning when the graphical model converges locally to a tree limit. A sufficient condition for the existence of such limits is the regime of *correlation decay*², which refers to the property that there are no long-range correlations in the model [29, 38, 51]. This regime is also known as the *uniqueness regime* since under such an assumption, the marginal distribution at a node is asymptotically independent of the configuration of a growing boundary.

We tailor the definition of correlation decay to node neighborhoods and provide the definition below. Given a graph $G = (W, E)$ and a graphical model $P_{\mathbf{X}_W|G}$ Markov on it, and any subset $A \subset W$, let $P_{\mathbf{X}_A|G}$ denote the marginal distribution of variables in A . For some subgraph $F \subset G$, let $P_{\mathbf{X}_A|F}$ denote the marginal distribution on A corresponding to a graphical model Markov on graph F instead of G (i.e., by setting the potentials of edges in $G \setminus F$ to zero). Let $\mathcal{N}[i; G] := \mathcal{N}(i; G) \cup i$ denote the closed neighborhood of node i in G . For any two sets $A_1, A_2 \subset W$, let $\text{dist}(A_1, A_2) := \min_{i \in A_1, j \in A_2} \text{dist}(i, j)$ denote the minimum graph distance³. Let $B_l(i)$ denote the set of nodes within graph distance l from node i and $\partial B_l(i)$ denote the boundary nodes, i.e., exactly at l from node i . Let $F_l(i; G) := G(B_l(i))$ denote the induced subgraph on $B_l(i)$. For any distributions P, Q , let $\|P - Q\|_1$ denote the ℓ_1 norm.

DEFINITION 1 (Correlation Decay). A graphical model $P_{\mathbf{X}_{W_m}|G}$ Markov on graph $G_m = (W_m, E_m)$ is said to exhibit correlation decay with a non-increasing rate function $\zeta_m(\cdot) > 0$ if for all $l, m \in \mathbb{N}$,

$$(5) \quad \|P_{\mathbf{X}_A|G_m} - P_{\mathbf{X}_A|F_l(i; G_m)}\|_1 \leq \zeta_m(\text{dist}(A, \partial B_l(i))), \quad \forall i \in W_m, A \subset B_l(i).$$

In words, the total variation distance⁴ between the marginal distribution of a set A of a graphical model Markov on G_m and the corresponding model Markov on subgraph $F_l(i; G_m)$ decays as a function of the graph distance to the boundary. This implies that for a class of functions $\zeta_m(\cdot)$, the effect of graph configuration beyond l hops from any node i has a decaying effect on the local marginal distributions.

For the class of Ising models in (4), the regime of correlation decay can be explicitly characterized, in terms of the maximum edge potential and the maximum degree of the graph, and this is studied in Section 3.2.

3. Method and Guarantees for Structure Estimation.

3.1. Overview of Algorithm. We now describe our algorithm, termed as **LocalCLGrouping**, for structure estimation of latent graphical models Markov on girth-constrained graphs. The algorithm leverages on the Chow-Liu grouping algorithm developed for latent tree models [19], described in Appendix A.1. The main intuition for learning a girth-constrained graph is based on reconstructing “local” parts of the graph which are acyclic and piecing them together. However, this approach has many challenges. First, it is not clear if the local acyclic pieces can be learnt efficiently since it requires the presence of an additive tree metric. This is addressed by considering models satisfying correlation decay (see Section 2.3). Second and a harder challenge involves merging the reconstructed local latent trees with provable guarantees due to the introduction of unlabeled latent

²Technically, correlation decay can be defined in multiple ways [38, p. 520] and the notion we use is the uniqueness or the weak spatial mixing condition.

³We distinguish between the terms *graph distance* and *information distances*. The former refers to the least number of hops on the graph, while the latter refers to the quantity in (6).

⁴Recall that the total variation distance between two probability distributions P, Q on the same alphabet is given by $\frac{1}{2}\|P - Q\|_1$.

Algorithm 1 LocalCLGrouping($\widehat{\mathbf{d}}^n(V), \Lambda, \tau, r$) for graph estimation using distance estimates $\widehat{\mathbf{d}}^n(V) := \{\widehat{d}(i, j)\}_{i, j \in V}$, confidence interval Λ , threshold τ and distance parameter r .

Input: Distance estimates between the observed nodes $\widehat{\mathbf{d}}^n(V) := \{\widehat{d}(i, j)\}_{i, j \in V}$, confidence parameter Λ , threshold τ and bound r on distances used for local reconstruction. Let $B_r(v; \widehat{\mathbf{d}}^n) := \{u : \widehat{d}^n(u, v) \leq r\}$ and $\text{MST}(A; \widehat{\mathbf{d}}^n)$ denotes the minimum spanning tree over $A \subset V$ based on edge weights $\widehat{\mathbf{d}}^n(A)$. Given a graph G , let $\text{Leaf}(G)$ denote the set of nodes with unit degree. Let $\mathcal{N}[i; G]$ denote the closed neighborhood of node i in graph G . $\text{RG}(\widehat{\mathbf{d}}^n(A), \Lambda, \tau)$ represents the recursive grouping method for building latent trees (see Appendix A.1) over the set of nodes A using distance estimates $\widehat{\mathbf{d}}^n(A)$ with confidence bound Λ and threshold τ for merging nodes.

```

for  $v \in V$  do
   $T_v \leftarrow \text{MST}(B_r(v); \widehat{\mathbf{d}}^n)$ .
end for
Initialize  $\widehat{G}, \widehat{G}_0 \leftarrow \cup_v T_v$ .
for  $v \in V \setminus \text{Leaf}(\widehat{G}_0)$  do
   $A \leftarrow \mathcal{N}[v; \widehat{G}]$ .
   $S \leftarrow \text{RG}(\widehat{\mathbf{d}}^n(A), \Lambda, \tau)$ .
   $\widehat{G}(A) \leftarrow S$  (Replace subgraph over  $A$  with  $S$  in  $\widehat{G}$ )
end for
Output  $\widehat{G}$ .

```

nodes in different pieces. We circumvent this challenge by leveraging on the Chow-Liu grouping algorithm [19] and merging the different pieces before introducing the latent nodes.

The algorithm is described in Algorithm 1. Let $\widehat{d}^n(i, j)$ denote the estimated distance between nodes i and j according to (43) using the empirical distribution $\widehat{P}_{\mathbf{x}_{i,j}}^n$ computed using n samples, i.e.,

$$(6) \quad \widehat{d}^n(i, j) := -\log |\det(\widehat{P}_{\mathbf{x}_{i,j}}^n)|, \quad \forall i, j \in V.$$

The set of distance estimates $\widehat{\mathbf{d}}^n(V) := \{\widehat{d}^n(i, j) : i, j \in V\}$ are input to the algorithm along with a parameter r . Recall that $B_r(i; \widehat{\mathbf{d}}^n(V)) := \{j : \widehat{d}^n(i, j) \leq r\}$. For each observed node $i \in V$, the set of nodes $B_r(i; \widehat{\mathbf{d}}^n(V))$ is considered, and the minimum spanning tree is constructed. The graph estimate \widehat{G}^n is initialized by taking the union of all the local minimum spanning trees. The latent nodes are now iteratively added by considering local neighborhoods of \widehat{G} and using any latent tree algorithm for reconstruction (e.g. [19, 39]). Note that the running time is polynomial (in the number of nodes) as long as polynomial time algorithms are employed for local latent tree reconstruction.

The proposed method is efficient for practical implementation due to the “divide and conquer” feature, i.e., the local latent tree building operations can be parallelized to obtain speedups. For real datasets, a tradeoff between model complexity and fidelity is typically enforced by optimizing scores such as the Bayesian information criterion (BIC) [47]. Such criteria can be easily enforced through a greedy local search in each iteration of our method, and this limits the number of hidden variables added by our method. In our experiments in Section 5, we found that this method is fast to implement on real and synthetic datasets.

We subsequently establish the correctness of the proposed method under a set of natural conditions. We require that the parameter r , which determines the set $B_r(i; \mathbf{d})$ for each node i , needs to be chosen as a function of the depth δ (i.e., distance from a hidden node to its closest observed nodes) and girth g of the graph. In practice, the parameter r provides flexibility in tuning the length of cycles added to the graph estimate. When r is large enough, we obtain a latent tree, while for small r , the graph estimate can contain many short cycles (and potentially many components). In experiments, we evaluate the performance of our method for different values of r . For more details, see Section 5.

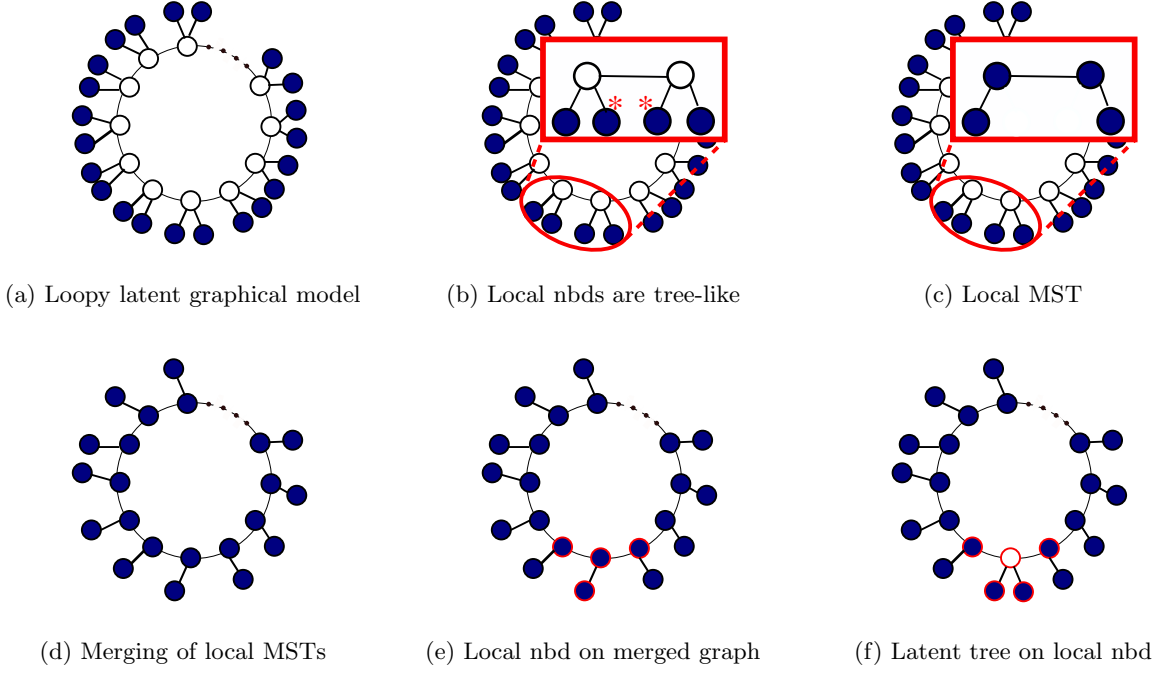


FIG 1. Various steps of LocalCLGrouping method on a simple cycle, where observed variables are shaded. See Section 3.1.1.

3.1.1. Toy Example with a Single Cycle. To demonstrate the steps of the above proposed method, consider the simple case of a single cycle of length g , where all the nodes on the cycle are hidden and each hidden node has two observed leaves, as shown in Fig.1a. When the cycle length g is sufficiently large, information distances on local neighborhoods are approximately additive, as depicted in Fig.1b. Moreover, in Fig.1b, let “*” denote the observed node closest to each hidden node (termed as its *surrogate*), in terms of information distance. The minimum spanning tree over the set of four nodes, which are zoomed in, corresponds to a chain shown in Fig.1c. Similarly, if in different local neighborhoods of observed nodes (based on a threshold on information distances), the surrogate relationships are similar (i.e., every hidden node has one of its children as its surrogate), then the local MSTs are simple chains, and their merging gives rise to graph G in Fig.1d. Now if a local neighborhood is selected on the merged graph G , as shown in Fig.1e, then we can discover the local latent tree structure based on information distances as shown in Fig.1f, since they are approximately additive. Similarly, when different neighborhoods on G are selected, local latent trees are discovered, and we recover the latent cycle graph in Fig.1a in the end.

3.2. Results for Ising Models. We first limit ourselves to providing asymptotic guarantees for the Ising model in (7), and then extend the results to non-asymptotic guarantees in general discrete distributions. Recall that the Ising model is a pairwise model over binary variables $\{-1, +1\}^p$ with the probability mass function (pmf) given by

$$(7) \quad P(\mathbf{x}_W) = \exp \left(\sum_{e \in E} \theta_{i,j} x_i x_j + \sum_{i \in V} \phi_i x_i - A(\boldsymbol{\theta}) \right).$$

3.2.1. *Conditions for Recovery in Ising Models.* We present a set of natural conditions on the graph structure and model parameters under which our proposed method succeeds in structure estimation.

- (A1) **Minimum Degree of Latent Nodes:** We require that all latent nodes have degree at least three, which is a natural assumption for identifiability of hidden variables. Otherwise, the latent nodes can be marginalized to obtain an equivalent representation of the observed statistics.
- (A2) **Distance Bounds:** Assume bounds on the edge potentials $\boldsymbol{\theta} := \{\theta_{i,j}\}$ of the Ising model:

$$(8) \quad \theta_{\min} \leq |\theta_{i,j}| \leq \theta_{\max}, \quad \forall (i,j) \in G.$$

Similarly assume bounded node potentials. We now define certain quantities which depend on the edge potential bounds. Given an Ising model P with edge potentials $\boldsymbol{\theta} = \{\theta_{i,j}\}$ and node potentials $\boldsymbol{\phi} = \{\phi_i\}$, consider its attractive counterpart \bar{P} with edge potentials $\bar{\boldsymbol{\theta}} := \{|\theta_{i,j}|\}$ and node potentials $\bar{\boldsymbol{\phi}} := \{|\phi_i|\}$. Let $\phi'_{\max} := \max_{i \in V} \text{atanh}(\bar{\mathbb{E}}(X_i))$, where $\bar{\mathbb{E}}$ is the expectation with respect to the distribution \bar{P} . Let $P(\mathbf{X}_{1,2}; \{\theta, \phi_1, \phi_2\})$ denote an Ising model on two nodes $\{1, 2\}$ with edge potential θ and node potentials $\{\phi_1, \phi_2\}$. Our learning guarantees depend on d_{\min} and d_{\max} satisfying

$$(9) \quad d_{\min} \geq -\log |\det P(\mathbf{X}_{1,2}; \{\theta_{\max}, \phi'_{\max}, \phi'_{\max}\})|,$$

$$(10) \quad d_{\max} \leq -\log |\det P(\mathbf{X}_{1,2}; \{\theta_{\min}, 0, 0\})|,$$

$$(11) \quad \eta := \frac{d_{\max}}{d_{\min}}.$$

- (A3) **Correlation Decay:** We assume correlation decay in the Ising model and require that

$$(12) \quad \alpha := \Delta_{\max} \tanh \theta_{\max} < 1, \quad \frac{\alpha^{g/2}}{\theta_{\min}^{\eta(\eta+1)+2}} = o(1),$$

where Δ_{\max} is the maximum node degree, g is the girth and $\theta_{\min}, \theta_{\max}$ are the minimum and maximum (absolute) edge potentials in the model.

- (A4) **Girth vs. Depth:** The depth δ characterizes how close the latent nodes are to observed nodes on graph G : for each hidden node $h \in H$, find a set of four observed nodes which form the shortest *quartet* with h as one of the middle nodes, and consider the largest graph distance in that quartet. The depth δ is the worst-case distance over all hidden nodes. We require the following tradeoff between the girth g and the depth δ :

$$(13) \quad \frac{g}{4} - \delta \eta (\eta + 1) = \omega(1),$$

Further, the parameter r in our algorithm is chosen as

$$(14) \quad r > \delta (\eta + 1) d_{\max} + \epsilon, \quad \text{for some } \epsilon > 0, \quad \frac{g}{4} d_{\min} - r = \omega(1).$$

- (A1) is a natural assumption on the minimum degree of the hidden nodes for identifiability. (A2) relates certain distance bounds to bounds on edge potentials. Intuitively, d_{\min} and d_{\max}

are bounds on information distances given by the local tree approximation of the loopy model, and its precise definition is given in (18). Note that $e^{-d_{\max}} = \Omega(\theta_{\min})$ and $e^{-d_{\min}} = O(\theta_{\max})$. (A3) uses bounds on the edge potentials to impose correlation decay on the model. It is natural that the sample requirement of any graph estimation algorithm depends on the “weakest” edge characterized by the minimum edge potential θ_{\min} . Further, the maximum edge potential θ_{\max} characterizes the presence/absence of long range correlations in the model. Intuitively, there is a tradeoff between the maximum degree Δ_{\max} and the maximum edge potential θ_{\max} of the model. Moreover, (A3) prescribes that the extent of correlation decay be strong enough (i.e., a small α and a large enough girth g) compared to the weakest edge in the model. Similar conditions have been imposed before for graphical model selection in the regime of correlation decay when there are no hidden variables [4]. (A4) provides the tradeoff between the girth g and the depth δ . Intuitively, the depth needs to be smaller than the girth to avoid encountering cycles during the process of graph reconstruction. Recall that the parameter r in our algorithm determines the neighborhood over which local MSTs are built in the first step. It is chosen such that it is roughly larger than the depth δ in order for all the hidden nodes to be discovered. The upper bound on r ensures that the distortion from an additive metric is not too large. The parameters for latent tree learning routines (such as confidence intervals for quartet tests) are chosen appropriately depending on d_{\min} and d_{\max} . See Section 3.3.

3.2.2. Guarantees for Ising Models. We now establish that the proposed method correctly estimates the graph structure of an Ising model in high dimensions. Recall that δ is the depth (distance from a hidden node to its closest observed nodes), θ_{\min} is the minimum (absolute) edge potential and $\eta = \frac{d_{\max}}{d_{\min}}$ is the ratio of distance bounds.

THEOREM 1 (Structural Consistency for Ising Models). *Under (A1)–(A4), the probability that the proposed method is structurally consistent tends to one, when the number of samples scales as*

$$(15) \quad n = \Omega \left(\theta_{\min}^{-\delta\eta(\eta+1)-2} \log p \right).$$

Proof: See Appendix B. □

Remarks: .

1. Thus, for learning Ising models on locally tree-like graphs, the sample complexity is dependent both on the minimum edge potential θ_{\min} and on the depth δ . Our method is efficient in high dimensions since the sample requirement is only logarithmic in the number of nodes p .
2. **Dependence on Maximum Degree:** For the correlation decay to hold (A3), we require $\theta_{\min} \leq \theta_{\max} = \Theta(1/\Delta_{\max})$. This implies that the sample complexity is at least $n = \Omega(\Delta_{\max}^{\delta\eta(\eta+1)+2} \log p)$.
3. **Comparison with Fully Observed Models:** In the special case when all the nodes are observed⁵ ($\delta = 1$), we strengthen the results for our method and establish that the

⁵In the trivial case, when all the nodes are observed and the graph is locally tree-like, our method reduces to thresholding of information distances at each node, and building local MSTs. The threshold can be chosen as $r = d_{\max} + \epsilon$, for some $\epsilon > 0$.

sample complexity is $n = \Omega(\theta_{\min}^{-2} \log p)$. This matches the best known sample complexity for learning fully observed Ising models [4, 31].

4. **Comparison with Learning Latent Trees:** Our method is an extension of latent tree methods for learning locally tree-like graphs. The sample complexity of our method matches the sample requirements for learning general latent tree models [19, 25, 39]. Thus, we establish that learning locally tree-like graphs is akin to learning latent trees in the regime of correlation decay.

3.3. Extension to General Discrete Models. We now extend the results to general discrete models and provide non-asymptotic sample requirement guarantees for success of our proposed method.

Local Tree Approximation: . We first define the notion of a local tree metric $\mathbf{d}_{\text{tree}}(V)$ computed by limiting the model to acyclic neighborhood subgraphs between the respective node pairs. Given a graph $G = (W, E)$, let $\text{tree}(i, j; G) := G(B_l(i) \cup B_l(j))$, for $l = \lfloor g/2 \rfloor - 1$, denote the induced subgraph on $B_l(i) \cup B_l(j)$, where g is the girth of the graph. Recall that $B_l(i; G)$ denotes the set of nodes within graph distance l from i in G . When $l < g/2 - 1$ no cycles are encountered and thus the induced subgraph $\text{tree}(i, j; G)$ is acyclic. Recall that $P_{\mathbf{x}_{i,j}|G}$ denotes the pairwise marginal distribution between i and j induced by the graphical model $P(\mathbf{x}_W)$ on graph $G(W)$. Let $P_{\mathbf{x}_{i,j}|\text{tree}(i,j)}$ denote the pairwise marginal distribution between i and j induced by considering only the subgraph $\text{tree}(i, j; G) \subset G$. Denote

$$(16) \quad d(i, j; \text{tree}) := -\log |\det P_{\mathbf{x}_{i,j}|\text{tree}(i,j)}|.$$

$$(17) \quad d(i, j; G) := -\log |\det P_{\mathbf{x}_{i,j}|G}|.$$

Denote $\mathbf{d}_{\text{tree}}(V) := \{d(i, j; \text{tree}) : i, j \in V\}$ and $\mathbf{d}(V) := \{d(i, j; G) : i, j \in V\}$. Note that for loopy graphs in general, $d(i, j; G)$ is different from $d(i, j; \text{tree})$. The learner has access only to the empirical versions $\hat{\mathbf{d}}(V)$ of the distances $\mathbf{d}(V)$, and thus the learner cannot estimate $\mathbf{d}_{\text{tree}}(V)$. However, we use $\mathbf{d}_{\text{tree}}(V)$ to characterize the performance of our algorithm, we list the relevant assumptions below.

3.3.1. Conditions on the Model Parameters.

- (B1) **Minimum Degree:** The minimum degree of any hidden node in the graph is $\text{Deg}_{\min}(H) \geq 3$.
- (B2) **Bounds on Local Tree Metric:** Given a graphical model $P_{\mathbf{x}_W|G}$ Markov on graph G , the pairwise marginal distribution $P_{\mathbf{x}_{i,j}|\text{tree}(i,j)}$ between any two neighbors $(i, j) \in G$ are non-singular⁶ and the distances $d(i, j; \text{tree}) := -\log |\det P_{\mathbf{x}_{i,j}|\text{tree}(i,j)}|$ satisfy

$$(18) \quad 0 < d_{\min} \leq d(i, j; \text{tree}) \leq d_{\max} < \infty, \quad \forall (i, j) \in G(W), \quad \eta := \frac{d_{\max}}{d_{\min}},$$

for suitable parameters d_{\min} and d_{\max} . We explicitly characterize d_{\min} and d_{\max} for Ising models.

⁶Note that $P_{\mathbf{x}_{i,j}|\text{tree}(i,j)}$ for $(i, j) \in G(W)$ is the probability distribution obtained by retaining only the node potentials ϕ_i and ϕ_j , and the edge potential $\theta_{i,j}$ and removing rest of the nodes. The distance is given by $d(i, j; \text{tree}) := -\log |\det P_{\mathbf{x}_{i,j}|\text{tree}(i,j)}|$.

(B3) **Regime of Correlation Decay:** The pairwise statistics of the graphical model converge locally to a tree limit according to Definition 1 with function $\zeta_m(\cdot)$ in (5) satisfying

$$(19) \quad 0 \leq \zeta \left(\frac{g}{2} - \frac{r}{d_{\min}} - 1 \right) < \frac{v}{|\mathcal{X}|^2},$$

where g is the girth, r is the distance bound parameter in **LocalCLGrouping**, $|\mathcal{X}|$ is the dimension of each variable, d_{\min}, d_{\max} are the distance bounds in (18) and

$$(20) \quad v := \min \left(d_{\min}, 0.5e^{-r}(e^{d_{\min}} - 1), e^{-0.5d_{\max}(\frac{r}{d_{\min}} + 2)}, \frac{g}{4}d_{\min} - r, r - d_{\max}\delta(\eta + 1) \right).$$

(B4) **Confidence Interval for Quartet Test:** The confidence interval in **Quartet**($\hat{\mathbf{d}}, \Lambda$) routine in Algorithm 2 is chosen as

$$(21) \quad \Lambda = \exp \left[-\frac{d_{\max}}{2} \left(\frac{r}{d_{\min}} + 2 \right) \right].$$

(B5) **Threshold for Merging Nodes:** The threshold τ in **RG**($\hat{\mathbf{d}}, \Lambda, \tau$) routine in Algorithm 3 is chosen as

$$(22) \quad \tau = \frac{d_{\min}}{2} - |\mathcal{X}|^2 \zeta_m \left(\frac{g}{2} - 1 \right) > 0,$$

where $|\mathcal{X}|$ is the dimension of the variable at each node and $\zeta_m(\cdot)$ is the correlation decay function according to (5).

(B1) is a natural assumption on the minimum degree of the hidden nodes for identifiability, which is also needed for latent trees. The assumption (B2) states that every edge has bounded distances under local tree approximations. Recall that in the special case of Ising models, this can be expressed via bounds on edge potentials. The assumption (B3) on correlation decay imposes a constraint on the rate function $\zeta(\cdot)$, in terms of the girth of the graph g , the distance threshold r used by the proposed method, the distance bounds d_{\min} and d_{\max} and depth δ . Recall that the depth δ characterizes how close the latent nodes are to observed nodes on graph G : for each hidden node $h \in H$, find a set of four observed nodes which form the shortest *quartet* with h as one of the middle nodes, and consider the largest graph distance in that quartet. The depth δ is the worst-case distance over all hidden nodes. (B3) implies that we require that the depth δ satisfies

$$(23) \quad \frac{g}{4}d_{\min} > \delta(\eta + 1)d_{\max}.$$

Similarly, (B3) imposes constraints on the parameter r used by the proposed algorithm for building local minimum spanning trees in the first step. (B3) implies that r needs to be chosen as

$$(24) \quad \delta(\eta + 1)d_{\max} < r < \frac{g}{4}d_{\min} - r.$$

Intuitively, the above constraint implies that r is relatively small compared to the girth of the graph and large enough for every hidden node to be discovered. This enables the proposed algorithm to correct reconstruct latent trees locally.

The confidence interval constraint in (B4) is based on the concentration bounds for the empirical distances. The threshold for merging nodes in (B5) ensures that spurious hidden nodes are not added. These conditions are inherited from latent tree algorithms.

3.4. *Guarantees for the Proposed Method.* We now establish that the **LocalCLGrouping** algorithm is structurally consistent under the above conditions.

THEOREM 2 (Structural Consistency of **LocalCLGrouping**). *Under assumptions (B1)-(B5), the **LocalCLGrouping** algorithm is structurally consistent with probability at least $1 - \kappa$, for any $\kappa > 0$, when the number of samples n available for learning satisfies*

$$(25) \quad n > \frac{2|\mathcal{X}|^2}{(v - |\mathcal{X}|^2 \zeta_m(\frac{g}{2} - \frac{r}{d_{\min}} - 1))^2} \left(4 \log p + |\mathcal{X}| \log 2 - \log \frac{\kappa}{7} \right),$$

where v is given by (20).

Remarks:

1. Thus, we provide PAC guarantees for reconstructing latent graphical models on girth-constrained graphs. The conditions required for success are mild and transparent, and along the lines of the conditions required for learning latent tree models. The conditions imposed on the girth of the graph are relatively mild. We require that the girth be roughly larger than the depth and that the correlation decay function $\zeta_m(\cdot)$ be sufficiently strong (B3). Thus, learning girth-constrained graphs is akin to learning latent tree models (in terms of sample and computational complexities) under a wide range of conditions.
2. One notable additional condition required for learning girth-constrained graphs in contrast to latent trees is the requirement of correlation decay (B3). However, we note that this is only a sufficient condition, and not necessary for learnability. For instance, the result in [22] establishes that the pairwise statistics converges locally to a tree limit for all attractive Ising models with strictly positive node potentials, but without any additional constraints on the parameters. Our results and analysis hold in such scenarios since we only require local convergence to a tree metric.
3. The results above are applicable for discrete models but can be extended to Gaussian models using the notion of *walk-summability* in place of correlation decay according to (5) (see [3]) and the negative logarithm of the correlation coefficient as the distance metric (see [19]). The results can also be extended to more general linear models such as multivariate Gaussian model, Gaussian mixtures and so on, along the lines of [2].

Proof: The detailed proof is given in Appendix C. It consists of the following main steps:

1. We first prove correctness of **LocalCLGrouping** under the tree limit (i.e., distances $\mathbf{d}_{\text{tree}}(V) := \{d(i, j; \text{tree})\}_{i, j \in V}$) and then show sample-based consistency. The latter is based on concentration bounds, along the lines of analysis for latent tree models [25, 39], with an additional distortion introduced due to the presence of a loopy graph.
2. We now briefly describe the proof establishing the correctness of **LocalCLGrouping** algorithm under \mathbf{d}_{tree} in girth-constrained graphs. Intuitively, the distances $d(i, j; \text{tree})$ correspond to a tree metric when the graph distance $\text{dist}(i, j) < g/2 - 1$, where g is the girth. Since **LocalCLGrouping** infers latent trees only locally, it avoids running into cycles and thus correctly reconstructs the local latent trees. The initialization step in **LocalCLGrouping** corresponds to the correct merge of this local latent trees under the

assumptions on parameter r in (24) and the correctness of **LocalCLGrouping** is established. \square

3.4.1. Guarantees under Uniform Sampling. We have so far given guarantees for graph reconstruction, given an arbitrary set of observed nodes in the graph. We now specialize the results to the case when there is a uniform sampling of nodes and provide learning guarantees. This analysis provides intuitions on the relationship between the fraction of sampled nodes and the resulting learning performance.

Let $\mathcal{G}_{\text{Girth}}(m; g, \Delta_{\min}, \Delta_{\max})$ denote the ensemble of graphs on m nodes with girth at least g and minimum degree $\Delta_{\min} \geq 3$ and maximum degree Δ_{\max} . Let $\rho := \frac{p}{m}$ denote the uniform sampling probability for selecting observed nodes. We have the following result on the depth δ . Define a constant $\epsilon_0 > 0$ as

$$(26) \quad \epsilon_0 = -\frac{\log(4m\Delta_{\max}(1-\rho)^{(\Delta_{\min}-1)^{g/2}})}{\log m}.$$

LEMMA 1 (Depth Under Uniform Sampling). *Given uniform sampling probability of ρ , for any $\epsilon \leq \max(0, \epsilon_0)$,*

$$(27) \quad \delta < \frac{1}{\log(\Delta_{\min}-1)} \left(\log \left[\frac{\log(4m^{1+\epsilon}\Delta_{\max})}{|\log(1-\rho)|} \right] \right) \quad w.p. \geq 1 - m^{-\epsilon}.$$

Proof: The proof is by straightforward arguments on binomial random variables and the union bound. See Appendix C.4. \square

Remarks:

1. Assuming that the girth satisfies $g > 2\delta(1 + d_{\max}/d_{\min})$ w.h.p., when the sampling probability and the degrees are both constant, then

$$(28) \quad \rho = \Theta(1), \Delta_{\min}, \Delta_{\max} = O(1) \Rightarrow \delta = O(\log \log m) \Rightarrow n = \Omega(\text{poly}(\log m)), \text{ w.h.p.}$$

On the other hand, with vanishing sampling probability, for $\beta \in [0, 1)$, we have

$$(29) \quad \rho = \Theta(m^{\beta-1}), \Delta_{\min}, \Delta_{\max} = O(1) \Rightarrow \delta = O(\log m) \Rightarrow n = \Omega(\text{poly}(m)), \text{ w.h.p.}$$

2. Recall that for Ising models, the best-case sample complexity of **LocalCLGrouping** for structural consistency (when $\eta = 1$ and $\theta_{\min} = \theta_{\max} = \Theta(1/\Delta_{\max})$) scales as

$$n = \Omega(\Delta_{\max}^{2(\delta+1)} \log p).$$

Thus, under uniform sampling, the sample complexity required for consistency scales as

$$n = \Omega \left(\Delta_{\max}^2 \left(\frac{\log p}{|\log(1-\rho)|} \right)^{4 \frac{\log \Delta_{\max}}{\log(\Delta_{\min}-1)}} \log p \right).$$

For the special case when the graph is regular ($\Delta_{\min} = \Delta_{\max}$), this reduces to

$$(30) \quad n = \Omega(\Delta_{\max}^2 \rho^{-2} (\log p)^3).$$

4. Necessary Conditions for Graph Estimation. We have so far provided sufficient conditions for recovering latent graphical models Markov on girth-constrained graphs. We now provide necessary conditions on the number of samples required by any algorithm to reconstruct the graph. Let $\hat{G}_n : (\mathcal{X}^{|V|})^n \rightarrow \mathcal{G}_m$ denote any deterministic graph estimator using n i.i.d. samples from node set V and \mathcal{G}_m is the set of all possible graphs on m nodes. We first define the notion of the graph edit distance.

DEFINITION 2 (Edit Distance). *Let G, \hat{G} be two graphs⁷ with adjacency matrices $\mathbf{A}_G, \mathbf{A}_{\hat{G}}$, and let V be the set of labeled vertices in both the graphs (with identical labels). Then the edit distance between G, \hat{G} is defined as*

$$\text{dist}(\hat{G}, G; V) := \min_{\pi} \|\mathbf{A}_{\hat{G}} - \pi(\mathbf{A}_G)\|_1,$$

where π is any permutation on the unlabeled nodes while keeping the labeled nodes fixed.

In other words, the edit distance is the minimum number of entries that are different in $\mathbf{A}_{\hat{G}}$ and in any permutation of \mathbf{A}_G over the unlabeled nodes. In our context, the labeled nodes correspond to the observed nodes V while the unlabeled nodes correspond to latent nodes H . We now provide necessary conditions for graph reconstruction up to certain edit distance.

THEOREM 3 (Necessary Condition). *For any deterministic estimator $\hat{G}_m : (\mathcal{X}^{m^\beta})^n \mapsto \mathcal{G}_m$ based on n i.i.d. samples from m^β observed nodes $\beta \in [0, 1]$ of a latent graphical model Markov on graph $G_m \in \mathcal{G}_{\text{Girth}}(m; g, \Delta_{\min}, \Delta_{\max})$ on m nodes with girth g , minimum degree Δ_{\min} and maximum degree Δ_{\max} , for all $\epsilon > 0$, we have*

$$(31) \quad \mathbb{P}[\text{dist}(\hat{G}_m, G_m; V) > \epsilon m] \geq 1 - \frac{|\mathcal{X}|^{nm^\beta} m^{(2\epsilon+1)m} 3^{\epsilon m}}{m^{0.5\Delta_{\min}m} (m - g\Delta_{\max}^g)^{0.5\Delta_{\min}m}},$$

under any sampling process used to choose the observed nodes.

Proof: The proof is based on counting arguments. See Section D for details. \square

Remarks:

1. The above result states that roughly

$$(32) \quad n = \Omega(\Delta_{\min} m^{1-\beta} \log m) = \Omega\left(\frac{\Delta_{\min}}{\rho} \log p\right)$$

samples are required for structural consistency. Thus, when $\beta = 1$ (constant fraction of observed nodes), logarithmic number of samples are necessary while when $\beta < 1$ (vanishing fraction of observed nodes), polynomial number of samples are necessary for reconstruction. From (30), recall that for Ising models, under uniform sampling of observed nodes, the best-case sample complexity of **LocalCLGrouping** (for homogeneous models on regular graphs with the degree Δ and $\theta_{\min} = \theta_{\max} = \Theta(1/\Delta)$) scales as

$$n = \Omega(\Delta^2 \rho^{-2} (\log p)^3),$$

and thus, nearly matches the lower bound on sample complexity in (32).

⁷We consider inexact graph matching where the unlabeled nodes can be unmatched. This is done by adding required number of isolated unlabeled nodes in the other graph, and considering the modified adjacency matrices [13].

5. Experiments. In this section we present experimental results on real and synthetic data. We evaluate performance in terms of perplexity, predictive perplexity and topic coherence, used frequently in topic modeling. In addition, we also study tradeoff between model complexity and data fitting through the Bayesian information criterion (BIC) [47]. Experiments are conducted using the 20 newsgroup data set, monthly stock returns from the S&P 100 companies and synthetic data. The datasets, software code and results are available at <http://newport.eecs.uci.edu/anandkumar>.

5.1. Experimental Setup.

Synthetic data: . We also generate samples from a known latent graphical model, shown previously in Fig.1, with a fixed depth $\delta = 1$, a fixed latent node degree $\Delta = 4$, and different girths $g = 10, 20, 30, \dots, 100$. The node potentials are kept at zero, while the edge potentials are chosen randomly in the range $[0.05, 0.2]$. This ensures that the model remains in the regime of correlation decay since the critical potential $\theta^* = \text{atanh}(\Delta^{-1}) = 0.2554 > 0.2$.

Newsgroup data: . We employ latent graphical models for topic modeling, i.e., modeling the relationships between various words co-occurring in documents. Each hidden variable in the model can be thought of as representing a topic, and topics and words in a document are drawn jointly from the graphical model. For a latent tree graphical model, topics and words are constrained to form a tree, while loopy models relax this assumption. We consider 16,242 binary samples of 100 keywords selected from the 20 newsgroup data. Each binary sample indicates the appearance of the given words in each posting. These samples are divided in to two equal groups, training and test sets for learning and testing purposes.

S&P data: . We also employ latent graphical models for financial modeling and in particular, for estimating the dependencies between the stock trends of different companies. The data set consists of monthly stock returns of 84 companies⁸ listed in S&P 100 index from 1990 to 2007. Experiments with this dataset allows us to demonstrate the performance of our algorithm on data using a Gaussian graphical model.

Methods: . We consider a regularized variant of the method proposed earlier for latent graphical model selection. Here, in every iteration, the decision to add hidden variables to a local neighborhood is based on the improvement of the overall BIC score. This allows us to tradeoff model complexity and data fitting. In addition, we obtain better generalization by avoiding overfitting. Note that our proposed method only deals with structure estimation and we use expectation maximization (EM) for parameter estimation. For the newsgroup data we compare the proposed method with the LDA model⁹.

Implementation: . The above method is implemented in MATLAB. We used the modules for LBP, made available with UGM¹⁰ package. The LDA models are learnt using the lda package¹¹.

⁸The 16 companies added after 1990 are dropped from the list of 100 companies listed in S&P 100 stock index for this analysis.

⁹Typically, LDA models the counts of different words in documents. Here, since we have binary data, we consider a binary LDA model where the observed variables are binary.

¹⁰These codes can be downloaded from <http://www.di.ens.fr/~mschmidt/Software/UGM.html>

¹¹<http://chasen.org/~daiti-m/dist/lda/>

Threshold selection r for our method: . Recall that the parameter r in our method controls the size of neighborhoods over which the local MSTs are constructed in the first step of our method. We earlier presented ranges of r , where recovery of the loopy structure is theoretically guaranteed (w.h.p). However, in practice, this range is unknown, since the model parameters are unknown to the learner, and also since there is no ground truth with respect to real datasets. Here, we present intuitive criterion for selecting the threshold based on the BIC score. First, it is important to search for the optimal threshold in the correct range, and this range is given by

$$(33) \quad r_{\max} := \max_{(i,j) \in V \times V} d(i,j), \quad r_{\min} := \max_{j \in V} \min_{i \in V} d(i,j),$$

if we disallow for disconnected components in the output graph. Note that if we choose $r \geq r_{\max}$, then the output is a latent tree. In our experiments, we choose one value above r_{\max} to find a reference tree model and compare it with other outcomes. For the 20 newsgroup dataset, we find that $r_{\min} = 2.3678$ and $r_{\max} = 12.2692$. Therefore, we choose $r \in \{3, 5, 7, 9, 11, 13\}$ for our experiments on newsgroup data. For the monthly stock returns data, $r_{\min} = 1.0337$ and $r_{\max} = 8.1172$, and we choose r from 1.1 to 8.2.

Performance Evaluation: . We evaluate performance based on the test perplexity [42] given by

$$(34) \quad \text{Perp-LL} := \exp \left[-\frac{1}{np} \sum_{k=1}^n \log P(\mathbf{x}^{\text{test}}(k)) \right],$$

where n is the number of test samples and p is the number of observed variables (i.e., words). Thus the perplexity is monotonically decreasing in the test likelihood and a lower perplexity indicates a better generalization performance. Along the lines of (34), we also evaluate the predictive perplexity [8]

$$(35) \quad \text{Pred-Perp-LL} := \exp \left[-\frac{1}{np} \sum_{k=1}^n \log P(\mathbf{x}_{\text{pred}}^{\text{test}}(k) | \mathbf{x}_{\text{obs}}^{\text{test}}(k)) \right],$$

where a subset of word occurrences $\mathbf{x}_{\text{obs}}^{\text{test}}$ is observed in test data and the performance of predicting the rest of words is evaluated. In our experiments, we randomly select half the words in test samples.

We also consider regularized versions of perplexity that capture tradeoff between model complexity and likelihood, given by

$$(36) \quad \text{Perp-BIC} := \exp \left[-\frac{1}{np} \text{BIC}(\mathbf{x}^{\text{test}}) \right],$$

where the BIC score [47] is defined as

$$(37) \quad \text{BIC}(\mathbf{x}^{\text{test}}) := \sum_{k=1}^n \log P(\mathbf{x}^{\text{test}}(k)) - 0.5(\text{df}) \log n,$$

where df is the degrees of freedom in the model. For a graphical model, we set $\text{df}^{\text{GM}} := m + |E|$, where m is the total number of variables (both observed and hidden) and $|E|$ is

the number of edges in the model. For the LDA model, we set $\text{df}^{\text{LDA}} := (p(m - p) - 1)$, where p is the number of observed variables (i.e., words) and $m - p$ is the number of hidden variables (i.e., topics). This is because a LDA model is parameterized by a $p \times (m - p)$ topic probability matrix and a $(m - p)$ -length Dirichlet prior. Thus, the BIC perplexity in (36) is monotonically decreasing in the BIC score, and a lower BIC perplexity indicates better tradeoff between model complexity and data fitting. However, the likelihood and BIC score in (34) and (36) are not tractable for exact evaluation in general graphical models since they involve the partition function. We employ loopy belief propagation (LBP) to evaluate them¹². Note that it is exact on a tree model and approximate for loopy models. Along the lines of predictive perplexity in (35), we also consider its regularized version

$$(38) \quad \text{Pred-Perp-BIC} := \exp \left[-\frac{1}{np} \text{BIC}(\mathbf{x}_{\text{pred}}^{\text{test}} | \mathbf{x}_{\text{obs}}^{\text{test}}) \right],$$

where the conditional BIC score is given by

$$(39) \quad \text{BIC}(\mathbf{x}_{\text{pred}}^{\text{test}} | \mathbf{x}_{\text{obs}}^{\text{test}}) := \sum_{k=1}^n \log P(\mathbf{x}_{\text{pred}}^{\text{test}}(k) | \mathbf{x}_{\text{obs}}^{\text{test}}(k)) - 0.5(\text{df}) \log n,$$

In addition, we also evaluate topic coherence, frequently considered in topic modeling. It is based on the average pointwise mutual information (PMI) score

$$(40) \quad \overline{\text{PMI}} := \frac{1}{45|H|} \sum_{h \in H} \sum_{\substack{i, j \in \mathcal{A}(h) \\ i < j}} \text{PMI}(X_i; X_j), \quad \text{PMI}(X_i; X_j) := \log \frac{P(X_i = 1, X_j = 1)}{P(X_i = 1)P(X_j = 1)},$$

where the set $\mathcal{A}(h)$ represents the “top-10” words associated with topic $h \in H$. The number of such word pairs for each topic is $\binom{10}{2} = 45$, and is used for normalization. In [43], it is found that the PMI scores are a good measure of human evaluated topic coherence when it is computed using an external corpus. It is also observed that using a related external corpus gives a high PMI. Hence, in our experiments, we choose a corpus containing news articles from the NYT articles bag-of-words dataset. This dataset has a vocabulary of 102660 words from 300,000 separate articles [27]. For LDA models, the top 10 words for each topic are selected based on the topic probability vector. For latent graphical models, we use the criterion of information distances on the learnt model to select the 10 nearest words for each topic.

5.2. Experimental Results.

Results for Synthetic Data: Outcome of the proposed method on synthetic data with threshold levels closer to r_{max} , defined in (33), discover the latent nodes closer to the actual numbers in the underlying model, whereas lower thresholds introduce more cycles and hidden variables. This is intuitive and occurs with real datasets as well. The normalized BIC scores (normalized with respect to n and p) of the loopy graphs improve with the number of samples n , as shown in Figure 2b. This is expected since the data becomes less noisy with

¹²The likelihood is evaluated using $P(\mathbf{x}_V) = \frac{P(\mathbf{x}_{V \cup H})}{P(\mathbf{x}_H | \mathbf{x}_V)}$, where $P(\mathbf{x}_H | \mathbf{x}_V)$ and $P(\mathbf{x}_{V \cup H})$ are computed using LBP, which is exact for trees. The above expression holds for any configuration of hidden variables \mathbf{x}_H , however we use the most likely hidden state to avoid numerical issues.

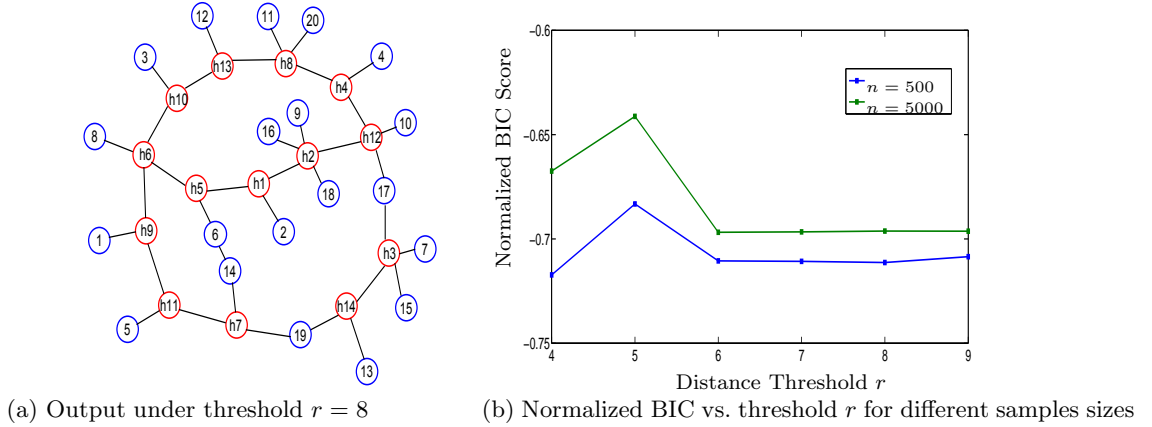


FIG 2. Results for synthetic data with girth $g = 10$ using the proposed method.

Method	r	Hidden	Edges	PMI	Perp-LL	Perp-BIC	Pred-Perp-LL	Pred-Perp-BIC
Proposed	3	55	265	0.2638	1.1533	1.1560	1.0695	1.0720
Proposed	5	39	293	0.4875	1.1567	1.1594	1.0424	1.0448
Proposed	7	32	183	0.4313	1.1498	1.1518	1.0664	1.0682
Proposed	9	24	129	0.6037	1.1543	1.1560	1.0780	1.0795
Proposed	11	26	125	0.4585	1.1555	1.1571	1.0787	1.0802
Proposed	13	24	123	0.4289	1.1560	1.1576	1.0788	1.0803
LDA	NA	10	NA	0.2921	1.1480	1.1544	1.1623	1.1656
LDA	NA	20	NA	0.1919	1.1348	1.1474	1.1572	1.1638
LDA	NA	30	NA	0.1653	1.1421	1.1612	1.1616	1.1715
LDA	NA	40	NA	0.1470	1.1494	1.1752	1.1634	1.1767

TABLE 1

Comparison of proposed method under different thresholds (r) with LDA under different number of topics (i.e., number of hidden variables) on 20 newsgroup data. For definition of perplexity and predictive perplexity based on test likelihood and BIC scores, and PMI, see (34), (35), (36), (38) and (40).

more samples. Figure 2b shows an overall improvement in the normalized BIC score with increasing number of samples n for different thresholds r . Figure 2b shows the variation of normalized BIC scores for graphs learnt using thresholds $r = 4$ to 9 with girth $g = 10$. We observe that the normalized BIC score decreases for the lowest threshold ($r = 4$), where the output graph shows a significant increase in latent nodes and edges, resulting in overfitting, and higher thresholds have better BIC. However, once the threshold results in a tree model, the BIC degrades since the cycles are no longer present.

Graph Structure for Newsgroup data: . We employ our method to learn the graph structures under different thresholds $r \in \{3, 5, 7, 9, 11, 13\}$ on newsgroup data, which controls the length of cycles. At $r = 13$ as shown in Fig 5, we obtain a latent tree and for $r \in \{3, 5, 7, 9\}$, we obtain loopy models. The first long cycle appears at $r = 9$ shown in Fig 4. At $r = 7$, we find a combination of short and long cycles. We find that models with cycles are more effective in discovering intuitive relationships. For instance, in the latent tree ($r = 13$), the link between “computer” and “software” is missing due to the tree constraint, but is discovered when $r \leq 9$. Moreover, we see that common words across different topics tend to connect the local subgraphs. For instance, the word “program” is used in the context of both space program

Topic 16	Topic 18	Topic 12	Topic 1	Topic 8
lunar	disk	card	god	software
moon	drive	video	jesus	pc
orbit	dos	windows	bible	computer
solar	memory	driver	christian	system
mission	windows	graphics	religion	dos
satellite	pc	dos	earth	windows
earth	software	version	question	disk
shuttle	scsi	ftp	fact	science
mars	computer	pc	jews	drive
space	system	disk	evidence	university

TABLE 2

Top 10 topic words from selected topics in loopy graphical model with threshold $r = 9$, the topic number corresponds to the labels of hidden variables in the loopy graph shown in Figure 4.

Topic 4	Topic 8	Topic 7	Topic 6	Topic 5
Space	windows	card	god	drive
nasa	files	graphics	world	states
insurance	dos	video	fact	research
earth	format	driver	christian	disk
moon	ftp	windows	jesus	university
orbit	program	computer	religion	mac
mission	software	pc	bible	scsi
launch	win	version	evidence	computer
gun	version	software	human	system
shuttle	pc	system	question	power

TABLE 3

Top 10 topic words corresponding to selected topics from the LDA model with 10 topics.

and computer programs. Similarly, the word “earth” is used both in the context of religion and space exploration.

Perplexity and Topic Coherence for Newsgroup Data: . In Table 1, we present results under our method and under LDA modeling on newsgroup data. For the LDA model, we vary the number of hidden variables (i.e., topics) as $\{10, 20, 30, 40\}$. In contrast, our method is designed to optimize for the number of hidden variables, and does not need this input. We note that our method is competitive in terms of both predictive perplexity and topic coherence. We find that the topic coherence (i.e., PMI) for our method is optimal at $r = 9$, where the graph has a single long cycle and a few short cycles. Intuitively, this model is able to discover more relationships between words, which the latent tree ($r = 13$) is unable to do so. On the other hand, for $r < 9$, topic coherence is degraded which suggests that adding too many cycles is counterproductive. However, the model at $r = 5$ performs better in terms of predictive perplexity indicating that it is able to use evidence from more observed words for prediction on test data. Moreover, all of our latent graphical models outperform the LDA models in terms of predictive perplexity. The top 10 topic words for selected topics are given for our method at ($r = 9$) and for the LDA model (with 10 topics) are given in Table 2 and Table 3.

Graph Structure for Stock Market Data: . The outcome of applying the proposed algorithm to stock market data is presented in Table 4. We observe that the number of edges and hidden variables remain fairly constant over a large range of thresholds. Specifically for $r \in [5.9, 6.7] \cup [6.8, 7.7]$, we obtain the same graph structure (for $r > r_{\max}$, we obtain a tree). Another

r	Hidden	Edges	Perp-LL	Perp-BIC
2.7	35	154	1.9498	2.0296
3.9	39	139	2.0200	2.0993
4.9	35	129	2.0210	2.0960
5	36	131	2.0169	2.0927
6.7	26	111	2.0344	2.1016
7.7	26	111	2.0353	2.1025
8.2	26	110	2.0405	2.1076

TABLE 4

Comparison of proposed method under different thresholds (r) on Stock data using the proposed method. For definition of perplexity based on test likelihood and BIC scores, see (34), and (36).

general trend observed is the improvement of the BIC score as the threshold decreases up till a certain point. The graphs learned using $r = 5, 7.7$ and 8.2 are shown in Fig.6, Fig.7, and Fig.8. Interesting connections between companies emerge. The latent tree structure in Fig.8 captures many key relationships. In particular, the S&P index node has a high degree since it captures the overall trend of various companies. Companies in similar sectors and divisions are grouped together. For instance, retail stores such as “Target”, “Walmart”, “CVS” and “Home Depot” are grouped together. However, additional relationships emerge as the threshold is decreased and cycles are added. We observe that the first cycle that is added connects the various oil companies which suggests strong interdependencies and influence on the S&P100 index. In addition, more cycles emerge when the threshold is decreased further. For instance, in Figure 6, we find a cycle connecting aviation company “Boeing” with “Honeywell” which is in aviation industry, but also additionally is in chemical industry and connects to companies such as “Dow Chemicals”. Thus, as in newsgroup data, we find that companies in multiple categories lead to cycles in the underlying graph.

Edge Density vs. Threshold r : We now study the edge density (i.e, number of edges) in the initialization step of our method as a function of the threshold r for both newsgroup and stock data. Recall that the initialization step involves building a loopy graph on observed variables (and no hidden variables). The edge density in this step is indicative of number of cycles added to the ultimate latent model. We observe that the graphs become denser as r is reduced from r_{\max} . However, when r is very small, the number of edges decreases since the nodes have sparser neighborhoods. This trend is seen in both : Figures 3a and 3b show the variation for newsgroup and stock data. For the newsgroup data, the graph density peaks at $r = 5$, which also achieves the highest predictive perplexity (see Table 1). Thus, we see a direct relationship between the edge density and the corresponding predictive perplexity in the learnt model. Intuitively, this is because as the number of edges increases, prediction at any node involves more evidence. However, as the threshold r is reduced further, graphs become less denser, and there is also a corresponding degradation in the predictive perplexity.

The above experiments confirm the effectiveness of our approach for discovering hidden topics, and are in line with the theoretical guarantees established earlier in the paper. Our analysis reveals that a large class of loopy graphical models with latent variables can be learnt efficiently in different domains.

6. Conclusion. In this paper, we considered latent graphical models Markov on girth-constrained graphs and proposed a novel approach for structure estimation. We established the correctness of the method when the model is in the regime of correlation decay and

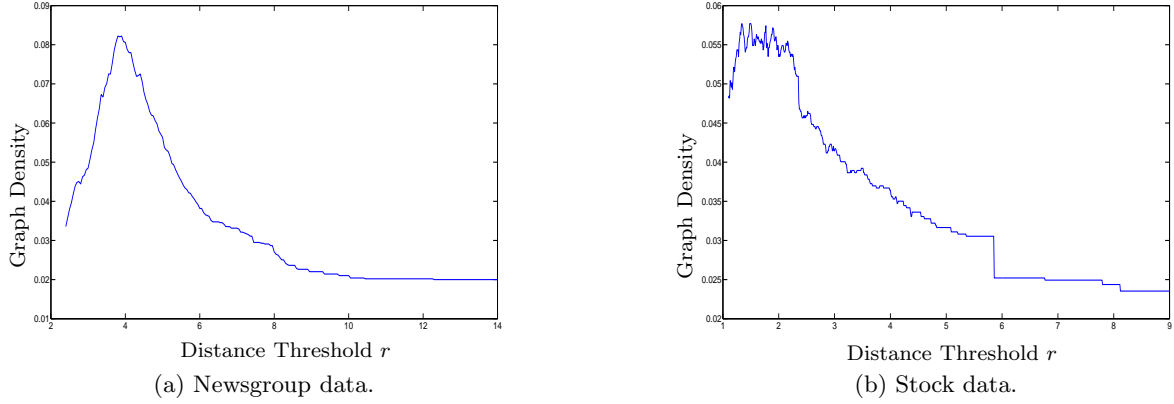


FIG 3. Variation of edge density of graphs at the initialization stage of LocalCLGrouping vs. threshold r .

also derived PAC learning guarantees. We compared these guarantees with other methods for graphical model selection, where there are no latent variables. Our findings reveal that latent variables do not add much complexity to the learning process in certain models and regimes, even when the number of hidden variables is large. These findings push the realm of tractable latent models for learning.

Acknowledgement. The author is supported in part by the setup funds at UCI and the AFOSR Award FA9550-10-1-0310. The authors thank E. Mossel (Berkeley) for detailed discussions in the beginning regarding problem formulation, modeling and algorithmic approaches, and Padhraic Smyth (UCI) and David Newman (UCI) for evaluation measures for topic models. The authors also thank the editor Tony Cai (Wharton) and anonymous reviewers whose comments substantially improved the paper.

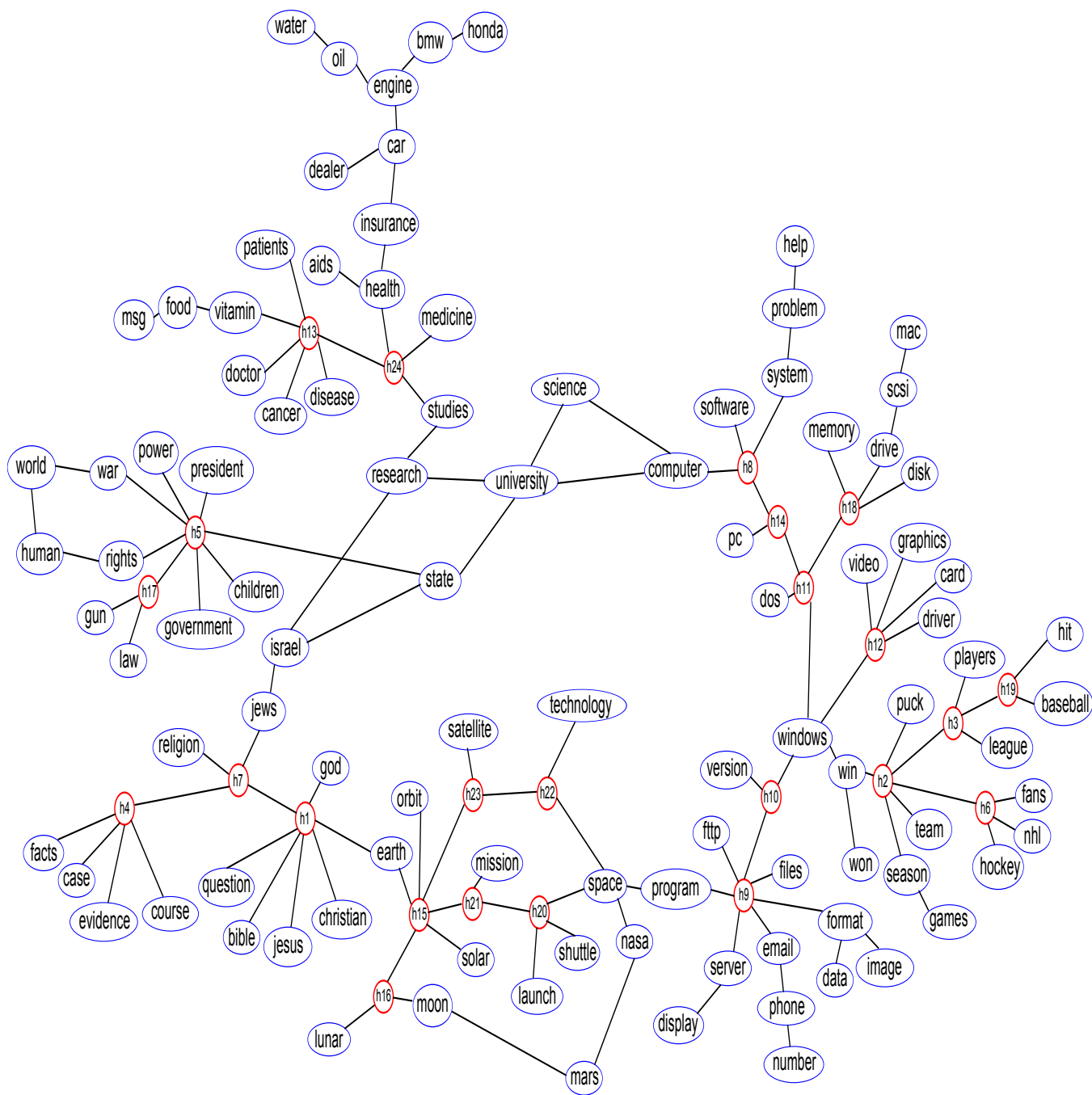


FIG 4. *Loopy Graph Learned using $r = 9$ with RegLocalCLGrouping on 20 newsgroup data.*



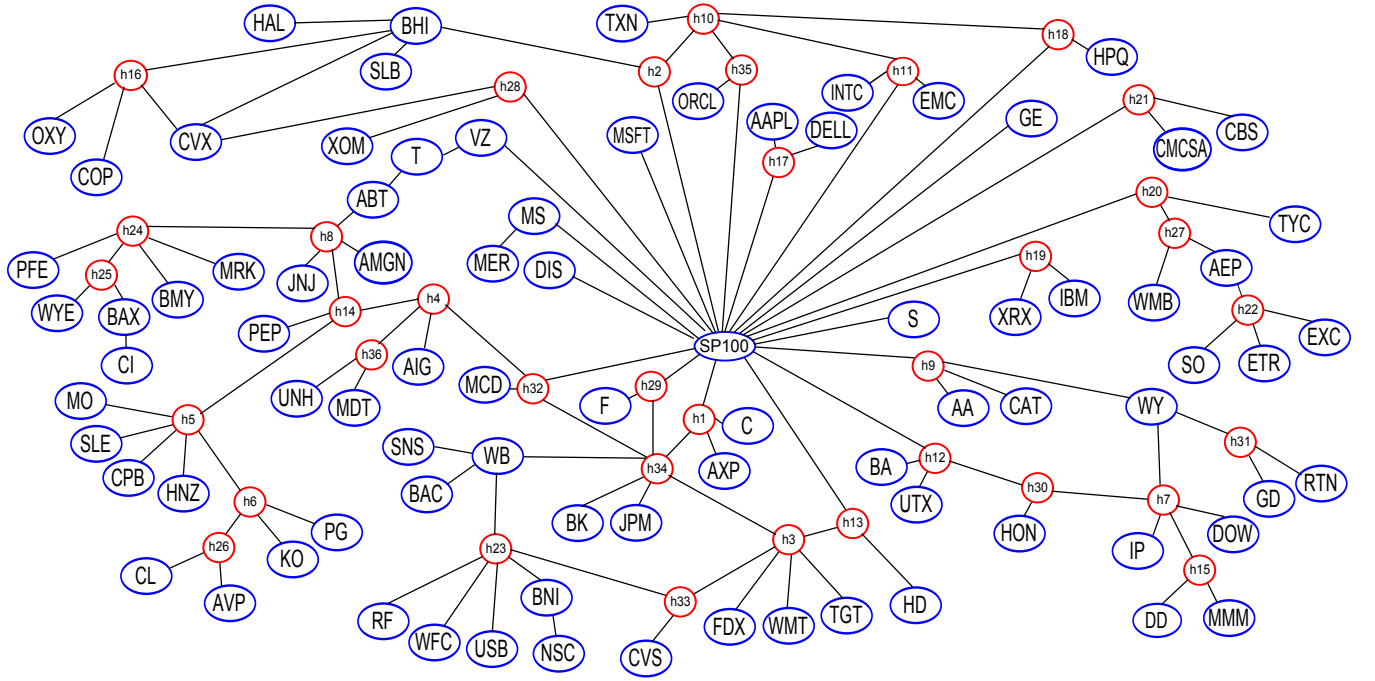


FIG 6. Loopy Graph Learned using $r = 5$ with LocalCLGrouping on S&P 100 monthly stock return data.

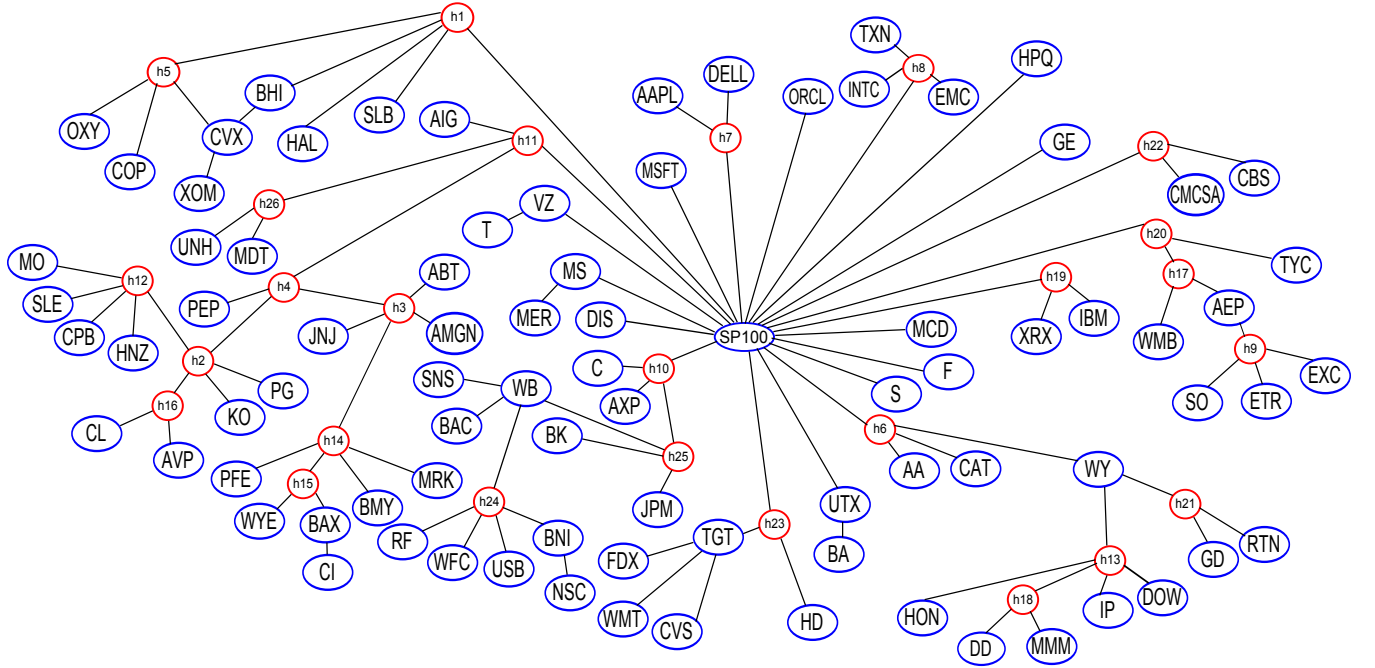


FIG 7. Loopy Graph Learned using $r = 7.7$ with LocalCLGrouping on S&P 100 monthly stock return data.

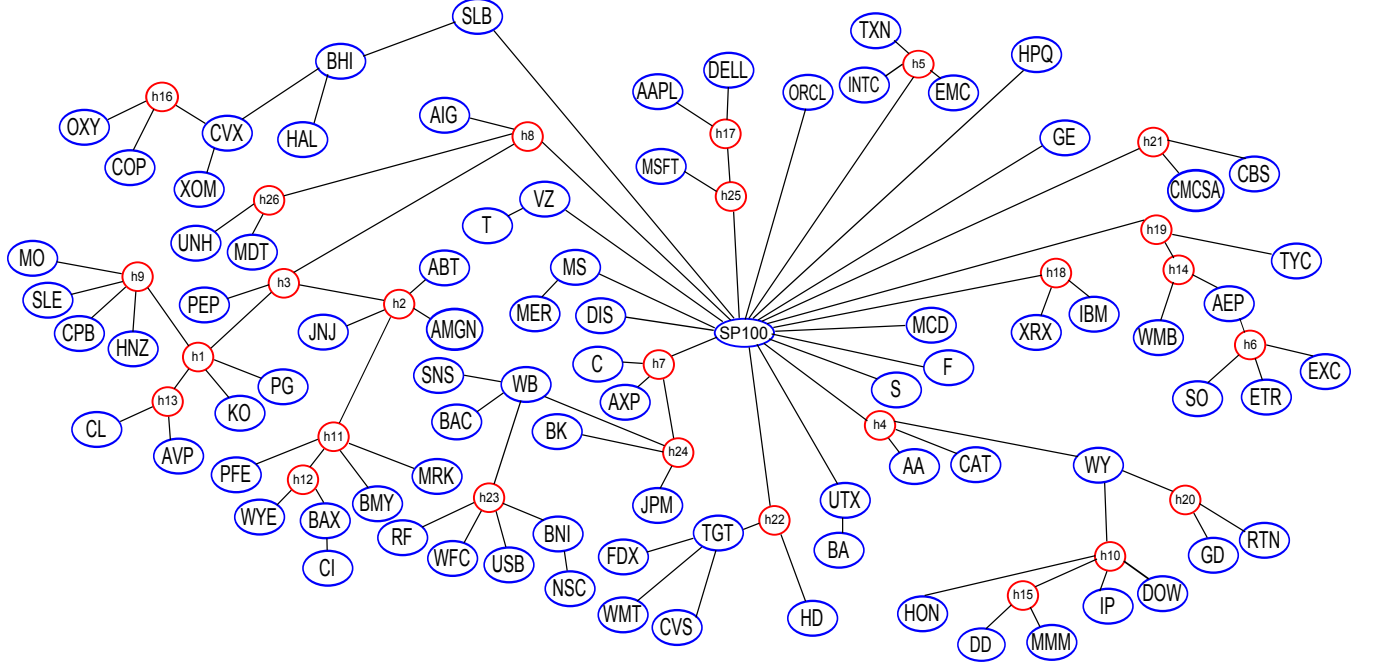


FIG 8. Tree Graph Learned using $r = 8.2$ with LocalCLGrouping on S&P 100 monthly stock return data.

APPENDIX A: BACKGROUND ON LATENT TREE MODELS

We first recap the results on latent tree models which will subsequently be extended to more general latent graphical models. It is well known that tree-structured graphical models Markov on a tree $T = (W, E)$ have a special form of factorization given by

$$(41) \quad P(\mathbf{x}_W) = \prod_{i \in W} P_{X_i}(x_i) \prod_{(i,j) \in T} \frac{P_{\mathbf{x}_{i,j}}(x_i, x_j)}{P_{X_i}(x_i)P_{X_j}(x_j)}$$

Comparing with general graphical models, we note that tree distributions are directly parameterized in terms of pairwise marginal distributions on the edges. Similarly, a Markov model can be described on a rooted directed tree \vec{T} with root $r \in W$, where the edges of \vec{T} are directed away from the root. Let $\text{Pa}(i)$ denote the (unique) parent of node $i \neq r$ and $P_{X_i|X_{\text{Pa}(i)}}$ denote the corresponding conditional distribution. The Markov model is given by

$$(42) \quad P(\mathbf{x}_W) = P_{X_r}(x_r) \prod_{i \in W, i \neq r} P_{X_i|X_{\text{Pa}(i)}}(x_i|x_{\text{Pa}(i)}).$$

A Markov model is said to be *non-singular* [40, 49] if (a) For all $e \in \vec{T}$, the conditional distributions satisfy $0 < |\det(P_{X_i|X_{\text{Pa}(i)}})| < 1$ and (b) For all $i \in V$, $P_{X_i}(x) > 0$ for all $x \in \mathcal{X}$.

A non-singular Markov model on an undirected tree T and its directed counterpart \vec{T} are equivalent [40, 49]. Note that non-singularity is equivalent to positivity (i.e., bounded potential functions) for Markov tree models. In particular, Ising models on trees with bounded node and edge potentials are non-singular. This is because under positivity, there is positive

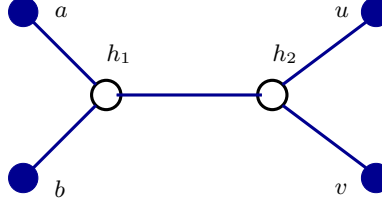


FIG 9. Quartet $Q(ab|uv)$. See (44).

probability for any global configuration of node states which implies that the conditional probability at a node given any of its neighbors cannot be degenerate.

Latent tree models or phylogenetic tree models are tree-structured graphical models in which a subset of nodes are hidden or latent. Our goal in this paper is to leverage on the techniques developed for learning latent tree models to analyze a more general class of latent graphical models.

A.1. Learning Latent Tree Models. Learning the structure of latent tree models is an extensively studied topic. A majority of structure learning methods (known as distance based methods) rely on the presence of an *additive tree metric* [23, 44]. The additive tree metric can be obtained by considering the pairwise marginal distributions of a tree structured graphical model. For instance, the work in [39] considers the following metric for discrete distributions satisfying the non-singular condition

$$(43) \quad d(i, j) := -\log |\det(P_{\mathbf{X}_{i,j}})|, \quad \forall i, j \in V.$$

By non-singularity assumption, we have that $|\det(P_{\mathbf{X}_{i,j}})| > 0$ for all $i, j \in W^2$. The distance metric further simplifies for some special distributions, e.g. for symmetric Ising models, it is given by the negative logarithm of the correlation between the node pair under consideration [48].

A.1.1. Quartet Based Methods. A popular class of learning methods are based on the construction of *quartets* (e.g., [12, 25, 39]), and various procedures to merge the inferred quartets. A quartet is a structure over four observed nodes, as shown in Fig.9. We now recap the classical quartet test operating on any additive tree metric [23, 44]. The path structure refers to the configuration of paths between the given nodes.

DEFINITION 3 (Quartet or Four-Point Condition on Trees). *Given an additive metric on a tree $[d(i, j)]_{i,j \in V}$, the tuple of four nodes $a, b, u, v \in V$ has the path structure in Fig.9 iff.*

$$(44) \quad d(a, b) + d(u, v) < \min(d(a, u) + d(b, v), d(b, u) + d(a, v)),$$

and the structure in Fig.9 is denoted by $Q(ab|uv)$.

It is well known that the set of all quartets uniquely characterize a latent tree. In [25], it was shown that a subset of quartets, termed as *representative quartets*, suffices to uniquely characterize a latent tree. The set of representative quartets consists of one quartet for each edge in the latent tree with shortest (graph) distances between the observed nodes.

Algorithm 2 $\text{Quartet}(\widehat{\mathbf{d}}^n(V), \Lambda)$ test using distance estimates $\widehat{\mathbf{d}}^n(V) := \{\widehat{d}(i, j)\}_{i, j \in V}$ and confidence interval Λ .

Input: Distance estimates between the observed nodes $\widehat{\mathbf{d}}^n(V) := \{\widehat{d}(i, j)\}_{i, j \in V}$ and confidence interval Λ .
Initialize set of quartets $\mathcal{Q}(V) \leftarrow \emptyset$.
for $\{i, j, i', j'\} \in V$ **do**
 if $(e^{-\widehat{d}(i, j)} - \Lambda)_+ (e^{-\widehat{d}(i', j')} - \Lambda)_+ > (e^{-\widehat{d}(i, j')} + \Lambda)_+ (e^{-\widehat{d}(i, j)} + \Lambda)_+$ **then**
 Declare Quartet: $\mathcal{Q}(V) \leftarrow \mathcal{Q}(V) \cup \{i, j, i', j'\}$.
 end if
 if No quartet declared for $\{i, j, i', j'\}$ **then**
 $\perp_{i, j, i', j'}$ (Declare null).
 end if
end for

Algorithm 3 $\text{RG}(\widehat{\mathbf{d}}^n(V), \Lambda, \tau)$ test using distance estimates $\widehat{\mathbf{d}}^n(V) := \{\widehat{d}(i, j)\}_{i, j \in V}$, confidence interval Λ and threshold τ for merging nodes.

Input: Distance estimates between the observed nodes $\widehat{\mathbf{d}}^n(V) := \{\widehat{d}(i, j)\}_{i, j \in V}$, confidence interval Λ and threshold τ . Let $\mathcal{C}(a)$ denote the children of node a .
Initialize $A \leftarrow V$, $\mathcal{C}(i) \leftarrow \{i\}$ for all $i \in V$ and $\mathcal{Q}(V) \leftarrow \text{Quartet}(\widehat{\mathbf{d}}^n(V), \Lambda)$.
while $A \neq \emptyset$ **do**
 if $\exists i, j \in A$ s.t. for each $a \in \mathcal{C}(i)$ and $b \in \mathcal{C}(j)$, $c, d \notin \mathcal{C}(i) \cup \mathcal{C}(j)$, $\{ac|bd, ad|bc\} \notin \mathcal{Q}(V)$, i.e., a, b are on same side of all such quartets in $\mathcal{Q}(V)$. **then**
 Declare i, j as siblings and introduce hidden node h as parent and $\mathcal{C}(h) \leftarrow \mathcal{C}(i) \cup \mathcal{C}(j)$.
 Remove i, j from A and add h to A .
 else
 Sibling relationships cannot be further inferred. Break.
 end if
end while
Form forest \widehat{T} based on sibling and child/parent relationships.
Merge edges in \widehat{T} of length less than τ and output \widehat{T} .

A.1.2. Recursive Grouping. We recap the recursive grouping $\text{RG}(\widehat{\mathbf{d}}^n(V), \Lambda, \tau)$ method proposed in [19] (and its refinement in [2]). The method is based on a robust¹³ quartet test $\text{Quartet}(\widehat{\mathbf{d}}^n, \Lambda)$ given in Algorithm 2. If the confidence bound is not met, a \perp result is declared. In the first iteration of RG, the algorithm searches for node pairs which occur on the same side of all the quartets, output by the quartet test $\text{Quartet}(\widehat{\mathbf{d}}^n, \Lambda)$ and declares them as siblings and introduces hidden variables. In later iterations of RG, sibling relationships between hidden variables are inferred through quartets involving their children. Finally, weak edges are merged and a tree (and more generally a forest) is output. We later use a modified version of recursive grouping method as a routine in our algorithm for estimating locally tree-like graphs. In the end, the neighboring nodes (at least one of which is hidden) are merged based on the threshold τ . See Section 3 for details.

A.1.3. Chow-Liu Grouping. An alternative method, known as *Chow-Liu grouping* (CLGrouping), was proposed in [19]. Although the theoretical results for CLGrouping are similar to earlier results (e.g. [25]), experiments on both synthetic and real data sets revealed significant improvement over earlier methods in terms of likelihood fitting and number of hidden variables added.

¹³Denote $(\cdot)_+ := \max(\cdot, 0)$.

Algorithm 4 CLGrouping($\widehat{\mathbf{d}}^n(V), \Lambda, \tau$) for graph estimation using distance estimates $\widehat{\mathbf{d}}^n(V) := \{\widehat{d}(i, j)\}_{i, j \in V}$, confidence interval Λ and threshold τ .

Input: Distance estimates between the observed nodes $\widehat{\mathbf{d}}^n(V) := \{\widehat{d}(i, j)\}_{i, j \in V}$, confidence interval Λ and threshold τ . Let $\text{MST}(V; \widehat{\mathbf{d}}^n)$ denote the minimum spanning tree over V according to the metric $\widehat{\mathbf{d}}^n(V)$. Given a tree T , let $\text{Leaf}(T)$ denote the set of leaves. Let $\mathcal{N}[i; T]$ denote the closed neighborhood of node i in tree T .

Initialize $\widehat{T} \leftarrow \text{MST}(V; \widehat{\mathbf{d}}^n)$.

for $v \in V \setminus \text{Leaf}(\widehat{T})$ **do**

$A \leftarrow \mathcal{N}[v; \widehat{T}]$.

$S \leftarrow \text{RG}(\widehat{\mathbf{d}}^n(A), \Lambda, \tau)$.

$\widehat{T}(A) \leftarrow S$ (Replace subtree over A with S in \widehat{T})

end for

Output \widehat{T} .

The CLGrouping method is summarized in Algorithm 4. The CLGrouping method always maintains a candidate tree structure and progressively adds more hidden nodes in local neighborhoods. The initial tree structure is the *minimum spanning tree* (MST) over the observed nodes with respect to the tree metric. The method then considers neighborhood sets on the MST and constructs local subtrees (using quartet based method or any other tree reconstruction algorithm). This local reconstruction property of CLGrouping makes it especially attractive for reconstructing girth-constrained graphs.

APPENDIX B: ANALYSIS OF ISING MODELS

For Ising models, the regime of correlation decay can be explicitly established. Recall that Δ_{\max} is the maximum degree of graph G and the maximum absolute edge potential is θ_{\max} .

LEMMA 2 (Correlation Decay in Ising Models). *The class of Ising models is in the regime of correlation decay when satisfies*

$$(45) \quad \alpha := \Delta_{\max} \tanh(\theta_{\max}) < 1.$$

The rate function $\zeta_m(\cdot)$ for correlation decay in (5) is given by

$$(46) \quad \zeta_m(l) = 2\alpha^l, \quad \forall l \in \mathbb{N}.$$

Moreover, for assumption (A3) to hold, it is sufficient that

$$(47) \quad \frac{\alpha^{g/2}}{\theta_{\min}^{\eta(\eta+1)+2}} = o(1),$$

where g is the girth, θ_{\min} is the minimum edge potential and $\eta := d_{\max}/d_{\min}$.

Proof: The above result on correlation decay is based on the concept of self-avoiding walk trees (SAW), which converts the conditional distributions of a general model to those on a tree model. See [4] for details. Regarding simplification of (B3) for Ising models to yield (A3), it is easy to see that $g/2 - r/d_{\min} - 1 = g/4 + \omega(1)$ from the constraint on r in (14). In Section B, it is shown that v is dominated by the second term, and its dependence on θ_{\min} is made explicit. \square

We now establish that the general distance bounds based on local tree approximation in (18) can be expressed in terms of potentials for the Ising model. Given an Ising model P with edge potentials $\boldsymbol{\theta} = \{\theta_{i,j}\}$ and node potentials $\boldsymbol{\phi} = \{\phi_i\}$, consider its attractive counterpart \bar{P} with edge potentials $\bar{\boldsymbol{\theta}} := \{|\theta_{i,j}|\}$ and node potentials $\bar{\boldsymbol{\phi}} := \{|\phi_i|\}$. Define the following quantity:

$$(48) \quad \phi'_{\max} := \max_{i \in V} \text{atanh}(\bar{\mathbb{E}}(X_i)),$$

where $\bar{\mathbb{E}}$ is the expectation with respect to the distribution \bar{P} . Finally let $P_{\mathbf{X}_{1,2};\{\theta,\phi_1,\phi_2\}}$ denote the Ising model on two nodes $\{1, 2\}$ with edge potential θ and node potentials ϕ_1 and ϕ_2 .

LEMMA 3 (Distance Bounds and Bounds on Edge Potentials). *For an Ising model, the distance bounds in (18) are related to the bounds on edge potentials in (8) as follows:*

$$(49) \quad d_{\min} \geq -\log |\det P_{\mathbf{X}_{1,2};\{\theta_{\max}, \phi'_{\max}, \phi'_{\max}\}}|,$$

$$(50) \quad d_{\max} \leq -\log |\det P_{\mathbf{X}_{1,2};\{\theta_{\min}, 0, 0\}}|.$$

The proof follows. We first consider marginal distribution on a tree model.

LEMMA 4 (Marginal Distribution on a Tree). *For an Ising model $\{\boldsymbol{\theta}, \boldsymbol{\phi}\}$ on a tree T with edge potentials $\boldsymbol{\theta} = \{\theta_{i,j}\}$ and node potentials $\boldsymbol{\phi} = \{\phi_i\}$, the marginal distribution between any two neighbors $(i, j) \in T$ is given by an Ising model $\{\theta_{i,j}, \phi'_i, \phi'_j\}$ with edge potential $\theta_{i,j}$ and node potentials ϕ'_i and ϕ'_j given by*

$$(51) \quad \phi'_i = \text{atanh}(\mathbb{E}[X_i; T_{-j}]), \quad \phi'_j = \text{atanh}(\mathbb{E}[X_j; T_{-i}]),$$

where T_{-j} denotes the tree with node j removed.

Proof: A more general version of this result is proven in [22, Lemma 4.1]. □

We also use the following result for attractive models ($\theta_{i,j} > 0, \forall (i, j) \in G$).

LEMMA 5 (Griffith's Second Inequality [30]). *For two attractive Ising models Markov on same graph $G = (W, E)$ with potentials $0 < \theta_{i,j} \leq \theta'_{i,j}$ for all $(i, j) \in G$, we have*

$$(52) \quad \mathbb{E} \left[\prod_{i \in U} X_i; \boldsymbol{\theta} \right] \leq \mathbb{E} \left[\prod_{i \in U} X_i; \boldsymbol{\theta}' \right], \quad \forall U \subset W.$$

In particular, this means that if the potentials of a model are increased, then the marginal expectation $\mathbb{E}[X_i]$ are also increased. This implies if some of the edge potentials are set to zero (meaning we take the model on a subgraph), $\mathbb{E}[X_i]$ is reduced.

Finally we note a simple expression for information distance in a symmetric Ising model (with zero node potentials) on two nodes.

FACT 1 (Symmetric Ising Model). *For a symmetric Ising model on two nodes $\{1, 2\}$ with edge potential θ and zero node potentials, we have*

$$(53) \quad d(1, 2) := -\log |\det P_{\mathbf{X}_{1,2};\{\theta, 0, 0\}}| = -\log |C_{1,2}| = -\log \tanh |\theta|,$$

where $C_{1,2} := \mathbb{E}[X_1 X_2]$ is the correlation between the two nodes.

Proof: For a symmetric model, we have $P(X_i = x) = 0.5$ for $i = 1, 2$ and $x \in \{-1, +1\}$. Similarly $P(X_1 = +|X_2 = -) = P(X_1 = -|X_2 = +)$. Using these facts, the distance $d(1, 2) = -\log |C_{1,2}|$. The correlation simplifies to

$$(54) \quad C_{1,2} := \frac{1}{2} (\mathbb{E}[X_1|X_2 = +] - \mathbb{E}[X_1|X_2 = -]) = \tanh \theta.$$

□

From the above fact, assuming Lemma 3 holds, d_{\max} for Ising models is given by

$$(55) \quad d_{\max} \leq -\log \tanh \theta_{\min}.$$

Proof of Lemma 3: For an Ising model $P_{\mathbf{X}_{1,2};\{\theta,\phi_1,\phi_2\}}$ on two nodes $\{1, 2\}$ with edge potential θ and node potentials ϕ_1, ϕ_2 , we have

$$(56) \quad \exp[-d(1, 2; \{\theta, \phi_1, \phi_2\})] = |\det P| = \frac{\sinh(2|\theta|)}{2(e^\theta \cosh(\phi_1 + \phi_2) + e^{-\theta} \cosh(\phi_1 - \phi_2))^2}.$$

Without loss of generality, consider an attractive model ($\theta > 0$). The above function is minimized with respect to $\{\phi_1, \phi_2\}$ when $\phi_1 = \phi_2 = 0$ since $\cosh(x) \geq \cosh(0) = 1$. Similarly it is maximized with respect to $\{\phi_1, \phi_2\}$ when $\phi_1 = \phi_2 = \phi'_{\max}$ for ($\theta > 0$). We subsequently establish that it is the maximum allowed node potential. When $\phi_1 = \phi_2 = 0$, we can show that $\exp[-d(1, 2)]$ is increasing in θ and thus, the minimum is attained when $\theta = \theta_{\min}$, and the maximum when $\theta = \theta_{\max}$.

From Lemma 4, the marginal distribution between two neighbors on a tree model is characterized. Only the node potentials at the two nodes are altered when the marginal distribution at the two nodes is considered. The (absolute) node potential at the two nodes is dominated by the attractive counterpart and cannot exceed ϕ'_{\max} in (51) from Griffith's property of attractive models in Lemma 5. □

Proof of Theorem 1: From Theorem 2, we have structural consistency when $n = \Omega(v^{-2} \log p)$, where v is given in (B3). We have $v = \min(-0.5e^{-r}(e^{d_{\min}} - 1), \exp[-0.5d_{\max}(r/d_{\min} + 2)])$. When r is chosen as $r = \delta(\eta + 1)d_{\max} + \epsilon$, for some $\epsilon > 0$, we have that $v = \exp[-0.5d_{\max}(r/d_{\min} + 2)]$, when $e^{-d_{\max}} < 1/3$. Using the fact that $e^{-d_{\max}} = \tanh \theta_{\min}$ from (55), we have that $\tanh \theta_{\min} < 1/3$ holds when the maximum degree $\Delta_{\max} > 3$ since the model is in the regime of correlation decay (B3). Since we require minimum degree of three for identifiability of hidden nodes (B1), this is satisfied, and we have the result.

For the special case when all the nodes are observed ($\delta = 1$), the sample complexity can be improved by selecting the parameter $r > d_{\max} + \epsilon$ for some $\epsilon > 0$, and only building local MSTs, and considering their union. In this case the sample complexity is given by $n = \Omega(e^{2r} \log p)$ which reduces to $n = \Omega(\theta_{\min}^{-2} \log p)$. □

APPENDIX C: STRUCTURAL CONSISTENCY OF LOCALCLGROUPING

We first establish that the LocalCLGrouping algorithm proposed in Section 3 recovers the unknown latent graph correctly when statistics corresponding to the tree limit are input. In Section C.2, we then establish that distances based on exact statistics converge locally to their tree limit. Finally, we consider sample-based analysis in Section C.3, and use standard concentration results, along the lines of [26, Section 6], and thereby proving Theorem 2.

C.1. Correctness of LocalCLGrouping under Local Tree Metric $\mathbf{d}_{\text{tree}}(V)$. Recall that $\mathbf{d}_{\text{tree}}(V) := \{d(i, j; \text{tree}) : i, j \in V\}$ is given by

$$d(i, j; \text{tree}) := -\log |\det P_{\mathbf{X}_{i,j} | \text{tree}(i,j)}|,$$

where $P_{\mathbf{X}_{i,j} | \text{tree}(i,j)}$ denotes the distribution at nodes i and j by limiting the model to the induced subgraph $\text{tree}(i, j)$. $\text{tree}(i, j) := G(B_l(i) \cup B_l(j))$ for $l = \lfloor g/2 \rfloor - 1$ and g is the girth of the graph. Since $\text{tree}(i, j)$ has no cycles, it immediately follows that $\mathbf{d}_{\text{tree}}(V)$ is a tree metric.

FACT 2 (Local Tree Metric). *The distances $\mathbf{d}_{\text{tree}}(V)$ form an additive tree metric.*

We now establish that the LocalCLGrouping algorithm proposed in Algorithm 1 outputs the correct graph under the assumptions on Theorem 2 when a local tree metric, computed using acyclic neighborhood subgraphs according to (16), $\mathbf{d}_{\text{tree}}(V)$ are input to the algorithm. Note that in practice, we only have access to empirical estimates $\hat{\mathbf{d}}^n(V)$ of the distances $\mathbf{d}(V)$, and not $\mathbf{d}_{\text{tree}}(V)$. In Section C.2, we establish the local convergence of $\mathbf{d}(V)$ to $\mathbf{d}_{\text{tree}}(V)$ under correlation decay.

C.1.1. Recap of CLGrouping for Learning Latent Trees. We first recap the result from [19, Lemma 8] that relates a latent tree model with the minimum spanning tree over the observed nodes according to a tree metric. Note that in this case, $\mathbf{d}(V)$ coincides with $\mathbf{d}_{\text{tree}}(V)$. For every node $i \in W$ in the latent tree T , define a mapping $\text{Sg} : W \mapsto V$, termed as *surrogate mapping*, as follows:

$$(57) \quad \text{Sg}(i; \mathbf{d}) := \arg \min_{j \in V} d(i, j; T), \quad \forall i \in W.$$

Thus, observed nodes V are their own surrogates while the hidden nodes H are mapped to the closest observed node according to metric $\mathbf{d}(V)$. See Fig.10 for an example.

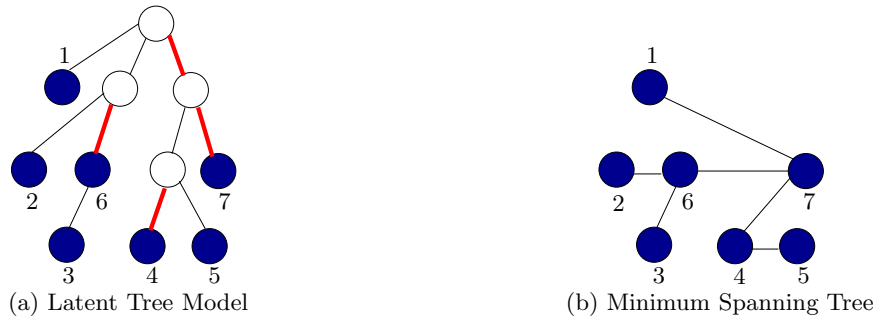


FIG 10. A latent tree model over $T = (W, E)$ and the corresponding minimum spanning tree $\text{MST}(V; \mathbf{d})$ over the observed nodes $V \subset W$. The observed nodes are shaded while the hidden nodes are unshaded. The thick lines in Fig.10a represent the edge between a hidden node and its surrogate. See Lemma 1.

PROPOSITION 1 (Relating Latent Tree and MST). *Given a latent tree $T = (W, E)$, set of observed nodes $V \subset W$ and a tree metric $\mathbf{d}(V)$, the minimum spanning tree $\text{MST}(V; \mathbf{d})$ over the observed nodes satisfies the following properties:*

1. The $\text{MST}(V; \mathbf{d})$ is obtained from the latent tree T by merging each hidden node $h \in H$ with its surrogate $\text{Sg}(h; \mathbf{d})$ and viceversa.
2. Let ξ denote the maximum graph distance between a hidden node and its surrogate in the latent tree T and let δ denote the depth of tree T . We have

$$(58) \quad \xi \leq \delta \frac{d_{\max;T}}{d_{\min;T}},$$

where $d_{\min;T}$ and $d_{\max;T}$ are bounds on the distance in T .

C.1.2. Union of Local MSTs under LocalCLGrouping. Using the results of **CLGrouping**, we establish properties of the union of local minimum spanning trees for girth-constrained graphs under correlation decay. To this end, consider the choice of parameter r in (24) and bounds $d_{\min;\text{tree}}$ and $d_{\max;\text{tree}}$. Also define

$$(59) \quad r' := \lfloor \frac{r}{d_{\max;\text{tree}}} \rfloor, \quad r'' := \lceil \frac{r}{d_{\min;\text{tree}}} \rceil.$$

Recall that $B_r(i; \mathbf{d}_{\text{tree}})$ denotes the set of observed nodes within distance r according to the metric $\mathbf{d}_{\text{tree}}(V)$. Let $B_{r'}(i; G)$ denote the set of nodes (including hidden nodes) within graph distance r' from node $i \in V$ on graph G . By definition, $B_{r'}(i; G) \subset B_r(i; \mathbf{d}_{\text{tree}}) \subset B_{r''}(i; G)$. In other words, the nodes in $B_r(i; \mathbf{d}_{\text{tree}})$ have graph distance at least r' and at most r'' . We have the following result.

LEMMA 6 (Properties of Union of Local MSTs under $\mathbf{d}_{\text{tree}}(V)$). *The graph formed by the union of local minimum spanning trees ($G' := \cup_{i \in V} \text{MST}(B_r(i; \mathbf{d}_{\text{tree}}))$) under LocalCLGrouping method using the distance metric $\mathbf{d}_{\text{tree}}(V)$, when the parameter r is chosen according to (24), satisfies the following properties:*

1. G' does not contain triangles.
2. G' is formed by contracting each hidden node $h \in H$ to its surrogate node $\text{Sg}(h; \mathbf{d}_{\text{tree}})$ (according to the distance metric (6)).

Proof: The first result is easy to see. We have that for each edge $(i, j) \in G'$, $d(i, j; \text{tree}) \leq r$ since the MSTs are formed on nodes within distance r . By contradiction, assume that a triangle exists between nodes $i, j, k \in V$ in G' . This implies that $d(i, j; \text{tree}), d(j, k; \text{tree}), d(k, i; \text{tree}) \leq r$. For a triangle to exist, we require another node $l \in V$ such that $d(j, l; \text{tree}), d(j, k; \text{tree}), d(k, l; \text{tree}) \leq r$. See Fig.11. Since the maximum graph distance between any two nodes i, j satisfying $d(i, j; \text{tree}) \leq r$ is r'' , we have that the maximum length of the cycle containing i, j, k, l is $4r''$. When $4r'' < g$ (which holds for r according to (24)), such a cycle cannot exist and such triangles cannot occur in G' .

For the second result, from Fact 2, the distances $\mathbf{d}_{\text{tree}}(B_{r''}(i; G))$ form a tree metric when $2r'' < g$, where g is the girth of the graph G , which holds for the choice of r in (24). This implies that Proposition 1 is applicable and the minimum spanning tree $\text{MST}(B_r(v); \mathbf{d})$ is formed as a result of contraction of hidden nodes to their surrogates. When the parameter ξ in (58) satisfies $\xi + \delta < r'$ (which is true under (23)), then every hidden node has a surrogate within some local neighborhood $B_r(v)$ and forms a quartet with its surrogate node. This

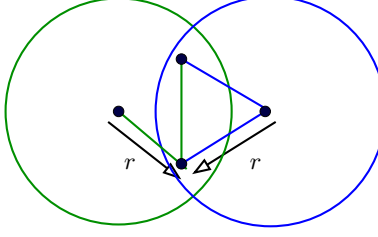


FIG 11. Condition for existence for triangles in $(G' := \cup_{i \in V} \text{MST}(B_r(i); \mathbf{d}_{\text{tree}}))$.

implies that every hidden node $h \in H$ contracts to its surrogate node in some local MST. \square

Proof of Theorem 2 under $\mathbf{d}_{\text{tree}}(V)$: We now show that the method **LocalCLGrouping** correctly recovers the graph G when tree-based distances $\mathbf{d}_{\text{tree}}(V)$ are input under the assumptions of Theorem 2. From Lemma 6, we have that in the graph formed from the union of local MSTs $(G' := \cup_{v \in V} \text{MST}(B_r(i); \mathbf{d}_{\text{tree}}))$, each hidden node is contracted to its surrogate node. The method **LocalCLGrouping** proceeds by reversing these contractions by considering neighborhoods on G' and constructing a local latent tree. Since there are no triangles in G' , the construction of local latent trees are independent. From the correctness of **CLGrouping** developed in [19], the local latent trees are correct since the distance metric converges locally to a tree metric. Thus, the correctness of **LocalCLGrouping** under $\mathbf{d}_{\text{tree}}(V)$ is proven. \square

Proof of Theorem 2 with samples: Combining Lemma 9, Lemma 9, and Lemma 12. \square

C.2. Local Convergence to a Tree Metric. We have so far analyzed the performance of **LocalCLGrouping** algorithm under tree-based distances $\mathbf{d}_{\text{tree}}(V)$. We now relate the distances $\mathbf{d}(V)$ computed using exact pairwise statistics with $\mathbf{d}_{\text{tree}}(V)$ under correlation decay according to (5). Let

$$d'_{\max}(l) := ld_{\max; \text{tree}} - \log(1 - e^{ld_{\max}} |\mathcal{X}| \zeta_m(g/2 - l - 1)),$$

where $d_{\max; \text{tree}}$ is the maximum $d(i, j; \text{tree})$ for any two neighbors i, j on graph G according to (18).

PROPOSITION 2 (Local Convergence to a Tree Metric). *When a discrete graphical model satisfies correlation decay with rate $\zeta_m(\cdot)$ according to (5), we have a.a.s., for nodes $i, j \in W$ with graph distance l in G and $l < g/2 - 1$,*

$$(60) \quad |\exp[-d(i, j; G)] - \exp[-d(i, j; \text{tree})]| \leq |\mathcal{X}| \zeta_m(g/2 - l - 1),$$

where g is the girth of the graph, and $|\mathcal{X}|$ is the cardinality of the random variable at each node. Additionally, we have

$$(61) \quad |d(i, j; G) - d(i, j; \text{tree})| \leq |\mathcal{X}| e^{d'_{\max}(l)} \zeta_m(g/2 - l - 1).$$

Proof: From the definition of correlation decay in (5), we have that

$$\|P_{\mathbf{X}_{i,j}|G} - P_{\mathbf{X}_{i,j}|\text{tree}(i,j)}\|_1 \leq \zeta_m(g/2 - l - 1),$$

since $\text{tree}(i, j; G) := G(B_{\lfloor g/2 \rfloor - 1}(i) \cup B_{\lfloor g/2 \rfloor - 1}(j))$ and $g/2 - l - 1$ is the distance from i and j to the nearest boundary.

From [7, Sec. 20], we have that for any $k \times k$ matrix A ,

$$(62) \quad |\det(A + E) - \det(A)| \leq k \max\{\|A\|_q, \|A + E\|_q\}^{k-1} \|E\|_q.$$

Thus, we have that

$$|\det(P_{\mathbf{X}_{i,j}|G}) - \det(P_{\mathbf{X}_{i,j}|\text{tree}(i,j)})| \leq |\mathcal{X}| \|P_{\mathbf{X}_{i,j}|G} - P_{\mathbf{X}_{i,j}|\text{tree}(i,j)}\|_1 \leq |\mathcal{X}| \zeta_m(g/2 - l - 1).$$

From Lipschitz continuity, we have that

$$|d(i, j; G) - d(i, j; \text{tree})| \leq e^{d'_{\max}(l)} |\det(P_{\mathbf{X}_{i,j}|G}) - \det(P_{\mathbf{X}_{i,j}|\text{tree}(i,j)})|,$$

Let $d_{\max;G}(l)$ be the maximum $d(i, j; G)$ for any two nodes i, j at graph distance l , and similarly for $d_{\max;\text{tree}}(l)$. Since no cycles are encountered in $\text{tree}(i, j)$, $d(i, j; \text{tree})$ is a tree metric and thus $d_{\max;\text{tree}}(l) = l d_{\max;\text{tree}}(1)$. For $d_{\max;G}(l)$, we note that

$$e^{-d_{\max;G}(l)} \geq e^{-l d_{\max;\text{tree}}(1)} - |\mathcal{X}| \zeta_m(g/2 - l - 1).$$

□

Remark: When

$$(63) \quad e^{l d_{\max;\text{tree}}} |\mathcal{X}| \zeta_m(g/2 - l - 1) = o(1),$$

then

$$(64) \quad |d(i, j; G) - d(i, j; \text{tree})| \leq |\mathcal{X}| e^{l d_{\max;\text{tree}} + o(1)} \zeta_m(g/2 - l - 1) = o(1),$$

C.3. Sample-Based Analysis.

C.3.1. Concentration of Distance Estimates. We first derive the concentration bounds for distance estimates along the lines of from [26, 39]. Let $\hat{\mathbf{d}}^n(V)$ be the estimated distances using n samples according to (6). We first recap the following result on empirical distribution [50, Thm. 2.1].

PROPOSITION 3 (Guarantees for General Empirical Distribution). *The following is true for the empirical distribution \hat{P}^n , obtained using n i.i.d. samples from a discrete distribution P :*

$$(65) \quad \mathbb{P}[\|\hat{P}^n - P\|_1 > \epsilon] \leq 2^k \exp[-n\epsilon^2/2],$$

where k is the dimension.

Given a graph G , let the graph distance between two nodes i and j under consideration on graph G be l . Recall that $|\mathcal{X}|$ is the dimension of the variable at each node.

LEMMA 7 (Concentration of Empirical Distances). *For empirical distance between node i and j at graph distance l , computed according to (6) using n samples, we have the following result:*

$$(66) \quad \mathbb{P} \left[\left| \exp[-\widehat{d}(i, j; G)] - \exp[-d(i, j; G)] \right| > \epsilon \right] \leq 2^{|\mathcal{X}|} \exp \left[-\frac{n\epsilon^2}{2|\mathcal{X}|^2} \right].$$

When $\epsilon > |\mathcal{X}|^2 \zeta_m(g/2 - l - 1)$ and $l < g/2 - 1$, we additionally have that

$$(67) \quad \mathbb{P} \left[\left| \exp[-\widehat{d}(i, j; G)] - \exp[-d(i, j; \text{tree})] \right| > \epsilon \right] \leq 2^{|\mathcal{X}|} \exp \left[-\frac{n}{2|\mathcal{X}|^2} \left(\epsilon - |\mathcal{X}|^2 \zeta_m \left(\frac{g}{2} - l - 1 \right) \right)^2 \right].$$

Proof: Along the lines of Proposition 2, using [7, Sec. 20], we have,

$$(68) \quad \mathbb{P} \left[\left| \det(\widehat{P}_{\mathbf{X}_{i,j}|G}^n) - \det(P_{\mathbf{X}_{i,j}|G}) \right| > \epsilon \right] \leq 2^{|\mathcal{X}|} \exp \left[-\frac{n\epsilon^2}{2|\mathcal{X}|^2} \right],$$

and thus (66) holds since $d(i, j; G) := -\log |\det(P_{\mathbf{X}_{i,j}|G})|$. Using Proposition 2, we also have (67). \square

C.3.2. *Recap of Sample Analysis of CLGrouping for Learning Latent Trees.* We also recap the result [19] that the minimum spanning tree (MST) constructed over observed nodes under CLGrouping method is consistent when the underlying model is a latent tree.

Recall that $p := |V|$ is the number of observed nodes and n is the number of samples. Let η be the maximum graph distance (with respect to the latent tree T) between any two neighbors in $\text{MST}(V, \mathbf{d})$ and d_{\min}, d_{\max} are distance bounds on the edges of the latent tree T .

LEMMA 8 (Consistency of MST using CLGrouping for Latent Trees). *Given a latent tree $T = (W, E)$ and observed node set $V \subset W$, the MST constructed by CLGrouping method using empirical distances $\widehat{\mathbf{d}}^n(V)$ does not coincide with the true MST based on exact distances $\mathbf{d}(V)$ with probability*

$$\mathbb{P} \left[\text{MST}(V; \widehat{\mathbf{d}}^n) \neq \text{MST}(V; \mathbf{d}) \right] \leq 2^{|\mathcal{X}|+1} p^3 \exp \left[-\frac{n}{8|\mathcal{X}|^2} e^{-2\eta d_{\max}} (1 - e^{-d_{\min}})^2 \right].$$

Proof: From the property of the MST,

$$\begin{aligned} \mathbb{P} \left[\text{MST}(V; \widehat{\mathbf{d}}^n) \neq \text{MST}(V; \mathbf{d}) \right] &\stackrel{(a)}{=} \mathbb{P} \left[\bigcup_{\substack{(i,j) \in \text{MST}(V; \mathbf{d}) \\ (i,j) \in \text{Path}(u,v)}} \left(e^{-\widehat{d}(u,v)} > e^{-\widehat{d}(i,j)} \right) \right], \\ &\stackrel{(b)}{\leq} p^3 \max_{\substack{(i,j) \in \text{MST}(V; \mathbf{d}) \\ (i,j) \in \text{Path}(u,v)}} \mathbb{P} \left[\epsilon_{u,v} - \epsilon_{i,j} > e^{-d(i,j)} - e^{-d(u,v)} \right] \\ &\stackrel{(c)}{\leq} p^3 \max_{i,j,u,v \in V} \mathbb{P} \left[\epsilon_{u,v} - \epsilon_{i,j} > e^{-\eta d_{\max}} (1 - e^{-d_{\min}}) \right] \end{aligned}$$

$$\stackrel{(d)}{\leq} 2p^3 \max_{u,v \in V} \mathbb{P} \left[\epsilon_{u,v} > \frac{e^{-\eta d_{\max}}}{2} (1 - e^{-d_{\min}}) \right],$$

where $\epsilon_{u,v} := \exp[-\widehat{d}(u, v)] - \exp[-d(u, v)]$ and similarly for $\epsilon_{i,j}$. Equality (a) is due to the property of the MST, inequality (b) is the union bound, inequality (c) is obtained by applying bounds on $d(i, j)$ and $d(u, v)$:

$$e^{-d(i,j)} - e^{-d(u,v)} > e^{-d(i,j)} (1 - e^{d(i,j)-d(u,v)}) > e^{-\eta d_{\max}} (1 - e^{-d_{\min}}),$$

since $d(i, j) \leq \eta d_{\max}$ for all $(i, j) \in \text{MST}(V; \mathbf{d})$ and $d(i, j) - d(u, v) \leq -d_{\min}$ for all $(i, j) \in \text{MST}(V; \mathbf{d})$ and $u, v \in V$ containing (i, j) on the path connecting them. Inequality (d) is obtained from the fact that $\epsilon_{u,v} - \epsilon_{i,j} \geq 2 \max(\epsilon_{u,v}, \epsilon_{i,j})$ and applying the union bound. The final result is from (68) in Proposition 7. \square

C.3.3. Sample Analysis of Union of MSTs under LocalCLGrouping. We now establish consistency under **LocalCLGrouping** algorithm using the above result and local convergence of the metric $\mathbf{d}(V)$ to tree-based metric $\mathbf{d}_{\text{tree}}(V)$, according to Proposition 2. Recall that $\widehat{\mathbf{d}}^n(V)$ denotes the estimates of the true distances $\mathbf{d}(V)$ according to graph G . Let $\mathbf{d}_{\text{tree}}(V)$ denote the distances by considering only acyclic neighborhood subgraphs, defined in Proposition 2. Given empirical distances $\widehat{\mathbf{d}}^n(V)$ and tree distances \mathbf{d}_{tree} and parameter r according to (24), for each $i \in V$, let $\widehat{A}_i := B_r(i; \widehat{\mathbf{d}}^n)$. Define $\mathcal{L} := \mathbb{N} \cap (r/d_{\min}, g/2)$.

LEMMA 9 (Union of Local MSTs under **LocalCLGrouping**: I). *Given a graphical model Markov on graph $G = (W, E)$ satisfying conditions of Theorem 2 with observed node set $V \subset W$, we have*

$$\begin{aligned} & \mathbb{P} \left[\bigcup_{i \in V} \text{MST}(\widehat{A}_i; \widehat{\mathbf{d}}^n) \neq \bigcup_{i \in V} \text{MST}(\widehat{A}_i; \mathbf{d}_{\text{tree}}) \right] \\ & \leq 2^{|\mathcal{X}|} p^3 \min_{l \in \mathcal{L}} \left(2p \exp \left[-\frac{n}{2|\mathcal{X}|^2} \left(0.5e^{-r}(e^{d_{\min}} - 1) - |\mathcal{X}|^2 \zeta_m\left(\frac{g}{2} - l - 1\right) \right)^2 \right] \right. \\ & \quad \left. + \exp \left[-\frac{n}{2|\mathcal{X}|^2} \left(l d_{\min} - r - |\mathcal{X}|^2 \zeta_m\left(\frac{g}{2} - l - 1\right) \right)^2 \right] \right). \end{aligned} \tag{69}$$

Remark: In the high-dimensional regime, where $p \rightarrow \infty$, the first term dominates. Since $\zeta_m(\cdot)$ is monotonically decreasing, we can choose $l = \lceil r/d_{\min} \rceil + 1$. Roughly, we require $n = \Omega(e^r)$ when the other parameters are bounded, for the error probability to decay.

Proof: Along the lines of Lemma 8, for each $k \in V$, we have

$$\begin{aligned} & \mathbb{P} \left[\text{MST}(\widehat{A}_k; \widehat{\mathbf{d}}^n) \neq \text{MST}(\widehat{A}_k; \mathbf{d}_{\text{tree}}) \right] \\ & = \mathbb{P} \left[\bigcup_{\substack{(i,j) \in \text{MST}(\widehat{A}_k; \mathbf{d}_{\text{tree}}) \\ (i,j) \in \text{Path}(u,v)}} \left(e^{-\widehat{d}(u,v)} > e^{-\widehat{d}(i,j)} \right) \right], \end{aligned}$$

$$\begin{aligned}
&\leq p^3 \max_{\substack{(i,j) \in \text{MST}(A_k; \mathbf{d}_{\text{tree}}) \\ (i,j) \in \text{Path}(u,v)}} \mathbb{P} [\epsilon_{u,v} - \epsilon_{i,j} > e^{-d(i,j;\text{tree})} - e^{-d(u,v;\text{tree})}] \\
&\stackrel{(a)}{\leq} p^3 \max_{i,j,u,v \in V} \mathbb{P} [\epsilon_{u,v} - \epsilon_{i,j} > (e^{-r} - \epsilon_{i,j})(1 - e^{-d_{\min}})] \\
&\stackrel{(b)}{\leq} 2p^3 \max_{u,v \in V} \mathbb{P} \left[\epsilon_{u,v} > \frac{e^{-r}}{2}(e^{d_{\min}} - 1) \right],
\end{aligned}$$

where $\epsilon_{u,v} := \exp[-\widehat{d}(u,v;G)] - \exp[-d(u,v;\text{tree})]$ and similarly for $\epsilon_{i,j}$. Inequality (a) is obtained by applying bounds on $d(i,j)$ and $d(u,v)$:

$$e^{-d(i,j)} - e^{-d(u,v)} > e^{-d(i,j)}(1 - e^{d(i,j)-d(u,v)}) > (e^{-r} - \epsilon_{i,j})(1 - e^{-d_{\min}}),$$

since $\widehat{d}(i,j) \leq r$ for all $(i,j) \in \text{MST}(V; \mathbf{d})$, $e^{-d(i,j)} > e^{-r} - \epsilon_{i,j}$ and $d(i,j) - d(u,v) \leq -d_{\min}$ for all $(i,j) \in \text{MST}(V; \mathbf{d})$ and $u,v \in V$ containing (i,j) on the path connecting them. Inequality (b) is obtained from the fact that $\epsilon_{u,v} - e^{-d_{\min}}\epsilon_{i,j} \geq 2 \max(\epsilon_{u,v}, e^{-d_{\min}}\epsilon_{i,j}) \geq 2e^{-d_{\min}} \max(\epsilon_{u,v}, \epsilon_{i,j})$ since $e^{-d_{\min}} < 1$.

Now define \widehat{l}_{\max} , the maximum graph distance between any two nodes in any \widehat{A}_k , i.e.,

$$\widehat{l}_{\max} := \max_k \left(\text{Diam}(\text{MST}(\widehat{A}_k)) \right),$$

where $\text{Diam}(\cdot)$ is the diameter, in terms of graph distance on G . From (67) in Lemma 7, conditioned on $\{\widehat{l}_{\max} = l\}$ and union bound on $k \in V$,

$$\begin{aligned}
&\mathbb{P} \left[\bigcup_{i \in V} \text{MST}(\widehat{A}_i; \widehat{\mathbf{d}}^n) \neq \bigcup_{i \in V} \text{MST}(\widehat{A}_i; \mathbf{d}_{\text{tree}}) \middle| \{\widehat{l}_{\max} = l\} \right] \\
(70) \quad &\leq 2^{|\mathcal{X}|+1} p^4 \exp \left[-\frac{n}{2|\mathcal{X}|^2} \left(0.5e^{-r}(e^{d_{\min}} - 1) - |\mathcal{X}|^2 \zeta_m \left(\frac{g}{2} - l - 1 \right) \right)^2 \right]
\end{aligned}$$

We now derive characterize the event that $\{\widehat{l}_{\max} = l\}$. Note that $\widehat{l}_{\max} \leq \widehat{d}'_{\max}/d_{\min}$, where

$$(71) \quad \widehat{d}'_{\max} := \max_{k \in V} \max_{i,j \in \widehat{A}_k} \widehat{d}(i,j).$$

Thus, we have

$$\begin{aligned}
\mathbb{P} [\widehat{l}_{\max} > l] &\leq \mathbb{P} \left[\bigcup_{\substack{k \in V \\ i,j \in \widehat{A}_k}} \widehat{d}(i,j) > ld_{\min} \right] \\
&\leq 2^{|\mathcal{X}|} p^3 \exp \left[-\frac{n}{2|\mathcal{X}|^2} \left(ld_{\min} - r - |\mathcal{X}|^2 \zeta_m \left(\frac{g}{2} - l - 1 \right) \right)^2 \right].
\end{aligned}$$

□

We now provide conditions when $\bigcup_{i \in V} \text{MST}(\widehat{A}_i; \mathbf{d}_{\text{tree}})$ coincides with $\bigcup_{i \in V} \text{MST}(A_i; \mathbf{d}_{\text{tree}})$, where $A_i := B_r(i; \mathbf{d}_{\text{tree}})$.

LEMMA 10 (Union of Local MSTs under LocalCLGrouping: II). *Given a graphical model Markov on graph $G = (W, E)$ satisfying conditions of Theorem 2 with observed node set $V \subset W$, we have*

$$\begin{aligned} & \mathbb{P} \left[\bigcup_{i \in V} \text{MST}(\hat{A}_i; \mathbf{d}_{\text{tree}}) \neq \bigcup_{i \in V} \text{MST}(A_i; \mathbf{d}_{\text{tree}}) \right] \\ & \leq 2^{|\mathcal{X}|} p^3 \left(\exp \left[-\frac{n}{2|\mathcal{X}|^2} \left(\frac{g d_{\min}}{4} - r - |\mathcal{X}|^2 \zeta_m \left(\frac{g}{2} - r/d_{\min} - 1 \right) \right)^2 \right] \right. \\ & \quad \left. + \exp \left[-\frac{n}{2|\mathcal{X}|^2} \left(r - \delta \left(\frac{d_{\max}}{d_{\min}} + 1 \right) d_{\max} - |\mathcal{X}|^2 \zeta_m \left(\frac{g}{2} - r/d_{\min} - 1 \right) \right)^2 \right] \right). \end{aligned}$$

Proof: Define

$$(72) \quad \hat{l}_{\max} := \max_k \left(\text{Diam}(\text{MST}(\hat{A}_k)) \right),$$

$$(73) \quad \hat{l}_{\min} := \min_k \left(\text{Diam}(\text{MST}(\hat{A}_k)) \right).$$

where $\text{Diam}(\cdot)$ is the diameter, in terms of graph distance on G . Conditioned on the event $\{\hat{l}_{\max} < \frac{g}{4}\} \cap \{\hat{l}_{\min} > \xi + \delta\}$, the graph satisfies the properties listed in Lemma 6 and thus,

$$(74) \quad \mathbb{P} \left[\bigcup_{i \in V} \text{MST}(\hat{A}_i; \mathbf{d}_{\text{tree}}) \neq \bigcup_{i \in V} \text{MST}(A_i; \mathbf{d}_{\text{tree}}) \mid \{\hat{l}_{\max} < \frac{g}{4}\} \cap \{\hat{l}_{\min} > \xi + \delta\} \right] = 0.$$

Moreover,

$$\begin{aligned} & \mathbb{P} \left[\{\hat{l}_{\max} > \frac{g}{4}\} \cup \{\hat{l}_{\min} < \xi + \delta\} \right] \\ & \leq 2^{|\mathcal{X}|} p^3 \left(\exp \left[-\frac{n}{2|\mathcal{X}|^2} \left(\frac{g d_{\min}}{4} - r - |\mathcal{X}|^2 \zeta_m \left(\frac{g}{2} - r/d_{\min} - 1 \right) \right)^2 \right] \right. \\ & \quad \left. + \exp \left[-\frac{n}{2|\mathcal{X}|^2} \left(r - (\xi + \delta) d_{\max} - |\mathcal{X}|^2 \zeta_m \left(\frac{g}{2} - r/d_{\min} - 1 \right) \right)^2 \right] \right), \end{aligned}$$

where $\xi := \delta d_{\max}/d_{\min}$ is the worst-case graph distance between a hidden node and its surrogate in G with respect to metric \mathbf{d}_{tree} . This is because the worst-case distance in a quartet containing a hidden node and its surrogate is $(\xi + \delta) d_{\max}$. When the empirical version of this distance exceeds r , then we have a bad event. \square

C.3.4. *Analysis of the Recursive Grouping.* Recall that for each $i \in V$, let $\hat{A}_i := B_r(i; \hat{\mathbf{d}}^n)$ and $A_i := B_r(i; \mathbf{d}_{\text{tree}})$. In LocalCLGrouping, the recursive grouping procedure is run on subsets of nodes in each \hat{A}_i . We first analyze the performance of quartet test.

LEMMA 11 (Analysis of Quartet Test). *Given distance estimates $\hat{\mathbf{d}}^n(\hat{A}_i)$ over observed nodes in \hat{A}_i , for each $i \in V$, $\text{Quartet}(\hat{\mathbf{d}}^n(\hat{A}_i), \Lambda)$ returns the correct set of quartets (and no null results) with probability at least*

$$\mathbb{P}[\cup_{i \in V} \{\text{Quartet}(\hat{\mathbf{d}}^n(\hat{A}_i), \Lambda) \neq \text{Quartet}(\mathbf{d}_{\text{tree}}(\hat{A}_i), \Lambda)\}]$$

$$(75) \quad \leq 2^{|\mathcal{X}|} p^3 \left(p \exp \left[-\frac{n}{2|\mathcal{X}|^2} \left(\exp[-(r/d_{\min} + 2)d_{\max}/2] - |\mathcal{X}|^2 \zeta_m(\frac{g}{2} - r/d_{\min} - 1) \right)^2 \right] \right. \\ \left. + \exp \left[-\frac{n}{2|\mathcal{X}|^2} \left(d_{\min} - |\mathcal{X}|^2 \zeta_m(\frac{g}{2} - r/d_{\min} - 1) \right)^2 \right] \right),$$

when Λ is chosen as

$$(76) \quad \Lambda = \exp[-(r/d_{\min} + 2)d_{\max}/2].$$

Proof: For each quartet $Q(v_1 v_2 | v_3 v_4)$ under metric $\mathbf{d}_{\text{tree}}(V)$ and $\mathcal{A} := \bigcup_{i=1}^4 v_i$, we have that

$$(77) \quad \mathbb{P}[\text{Quartet}(\hat{\mathbf{d}}^n(\mathcal{A}), \Lambda) \neq \text{Quartet}(\mathbf{d}_{\text{tree}}(\mathcal{A}), \Lambda) \mid \bigcap_{a,b \in \mathcal{A}} \{|\hat{d}^n(a,b) - d(a,b; \text{tree})| < \Lambda\}] = 0,$$

and the test $\text{Quartet}(\hat{\mathbf{d}}^n(\mathcal{A}), \Lambda)$ does not return null when $\Lambda < \exp[-\max_{a,b \in \mathcal{A}} d(a,b; \text{tree})/2]$. Considering all sets \hat{A}_i for $i \in V$, we require $\Lambda < \exp[-\hat{l}_{\max} d_{\max}/2]$ to not return null, where

$$\hat{l}_{\max} := \max_k \left(\text{Diam}(\text{MST}(\hat{A}_k)) \right).$$

From Lemma 7, choosing $\Lambda = \exp[-(l+1)d_{\max}/2]$ we that

$$\mathbb{P} \left[\bigcup_{i \in V} \{ \text{Quartet}(\hat{\mathbf{d}}^n(\hat{A}_i), \Lambda) \neq \text{Quartet}(\mathbf{d}_{\text{tree}}(\hat{A}_i), \Lambda) \mid \{\hat{l}_{\max} < l\} \} \right] \\ \leq 2^{|\mathcal{X}|} p^4 \exp \left[-\frac{n}{2|\mathcal{X}|^2} \left(\exp[-(l+1)d_{\max}/2] - |\mathcal{X}|^2 \zeta_m(\frac{g}{2} - l - 1) \right)^2 \right].$$

Along the lines of analysis in Lemma 10, we have that

$$\mathbb{P} \left[\{\hat{l}_{\max} > l\} \right] \leq 2^{|\mathcal{X}|} p^3 \exp \left[-\frac{n}{2|\mathcal{X}|^2} \left(l d_{\min} - r - |\mathcal{X}|^2 \zeta_m(\frac{g}{2} - l - 1) \right)^2 \right].$$

Choosing l as $r/d_{\min} + 1$, we have the result. \square

This yields the following result on recursive grouping.

LEMMA 12 (Results for Recursive Grouping). *The recursive grouping method $\text{RG}(\hat{\mathbf{d}}(\hat{A}_i), \Lambda, \tau)$ returns the same tree as $\text{RG}(\mathbf{d}_{\text{tree}}(\hat{A}_i), \Lambda, \tau)$ with probability*

$$(78) \quad \mathbb{P}[\bigcup_{i \in V} \{ \text{RG}(\hat{\mathbf{d}}^n(\hat{A}_i), \Lambda, \tau) \neq \text{RG}(\mathbf{d}_{\text{tree}}(\hat{A}_i), \Lambda, \tau) \}] \\ \leq 2^{|\mathcal{X}|} p^3 \left(p \exp \left[-\frac{n}{2|\mathcal{X}|^2} \left(\exp \left[-\frac{(\frac{r}{d_{\min}} + 2)d_{\max}}{2} \right] - |\mathcal{X}|^2 \zeta_m(\frac{g}{2} - \frac{r}{d_{\min}} - 1) \right)^2 \right] \right. \\ \left. + \exp \left[-\frac{n}{2|\mathcal{X}|^2} \left(d_{\min} - |\mathcal{X}|^2 \zeta_m(\frac{g}{2} - \frac{r}{d_{\min}} - 1) \right)^2 \right] \right. \\ \left. + \exp \left[-\frac{n}{2|\mathcal{X}|^2} \left(\frac{d_{\min}}{2} - |\mathcal{X}|^2 \zeta_m(\frac{g}{2} - 1) \right)^2 \right] \right),$$

when Λ is chosen as

$$(79) \quad \Lambda = \exp[-(r/d_{\min} + 2)d_{\max}/2].$$

and τ is chosen as

$$(80) \quad \tau = \frac{d_{\min}}{2} - |\mathcal{X}|^2 \zeta_m(g/2 - 1).$$

Proof: Along the lines of analysis in [2], given the correct set of quartets, the recursive grouping procedure returns the correct tree structure when the nodes are merged correctly with threshold τ . It is easy to see that this happens with probability

$$2^{|\mathcal{X}|} p^2 \exp \left[-\frac{n}{2|\mathcal{X}|^2} (d_{\min}/2 - |\mathcal{X}|^2 \zeta_m(g/2 - 1))^2 \right],$$

when the threshold is chosen as (80). □

Proof of Theorem 2: From Lemma 9, Lemma 10 and Lemma 12. □

C.4. Analysis Under Uniform Sampling. *Proof of Lemma 1:* Let $\mathcal{A}(e)$ denote the event that an hidden edge (with at least one hidden end point) has a representative quartet in which the end points are at most graph distance $l < g/2 - 1$. We have that

$$\mathbb{P}[\mathcal{A}^c(e)] \leq 4(1 - \rho)^{(\Delta_{\min} - 1)^{l-1}},$$

since there are at least $(\Delta_{\min} - 1)^{l-1}$ nodes in each of the four subtrees from which four observed nodes can be sampled and $\rho := p/m$ is the sampling probability. Taking the union bound, we have the probability that the depth δ is greater than $l < g/2 - 1$ as

$$\mathbb{P}[\delta > l] \leq 4m\Delta_{\max}(1 - \rho)^{(\Delta_{\min} - 1)^l}.$$

Thus, the result in (27) holds. □

APPENDIX D: NECESSARY CONDITIONS FOR GRAPH RECONSTRUCTION

Proof of Theorem 3: The proof is based on counting arguments along the lines of [11, Thm. 1]. For any deterministic estimator \hat{G}_m , let $\mathcal{R} := \hat{G}_m((\mathcal{X}^{m^\beta})^n)$ as the range of the estimator \hat{G}_m , when the number of observed nodes is $|V| = m^\beta$ for $\beta \in (0, 1]$. Thus, we have $|\mathcal{R}| = |\mathcal{X}|^{nm^\beta}$.

For any fixed graph F_m and set of labeled nodes V , denote the set of graphs within graph distance ϵm as

$$\mathcal{D}(F_m; \epsilon m) := \{G_m : \text{dist}(F_m, G_m; V) \leq \epsilon m\}.$$

We note that

$$|\mathcal{D}(F_m; \epsilon m)| \leq m! \binom{m^2}{\delta m} \leq m^{(2\epsilon+1)m} 3^{\epsilon m},$$

since we can permute the m vertices and change at most ϵm entries in the adjacency matrix \mathbf{A}_F and we use the bound that $\binom{N}{k} \leq \frac{N^k}{k!} \leq N^k 3^k$.

Let $\mathcal{S}(\hat{G}_m; \epsilon m)$ denote all the graphs which are within edit distance of ϵm of the graphs in range \mathcal{R} . We have that

$$|\mathcal{S}(\hat{G}_m; \epsilon m)| \leq |\mathcal{X}|^{nm^\beta} m^{(2\epsilon+1)m} 3^{\epsilon m}.$$

Along the lines of [11, Thm. 1], we have that the probability of error should satisfy

$$\mathbb{P}[\text{dist}(\hat{G}_m, G_m; V) > \epsilon m] \geq 1 - \frac{|\mathcal{S}(\hat{G}_m; \epsilon m)|}{|\mathcal{G}(m)|},$$

where $|\mathcal{G}(m)|$ is the number of graphs in the family under consideration.

From [4, Lemma 2], we have that for girth-constrained ensembles $\mathcal{G}_{\text{Girth}}(m; g, \Delta_{\min}, \Delta_{\max}, k)$ with girth g , minimum degree Δ_{\min} , maximum degree Δ_{\max} and number of edges k , we have

$$(81) \quad m^k (m - g\Delta_{\max}^g)^k \leq |\mathcal{G}_{\text{Girth}}(m; g, \Delta_{\min}, \Delta_{\max}, k)| \leq m^k (m - \Delta_{\min}^g)^k,$$

and we have the result. \square

References.

- [1] ANANDKUMAR, A., HASSIDIM, A. and KELNER, J. (2012). Topology Discovery of Sparse Random Graphs With Few Participants. *Accepted to J. of Random Structures and Algorithms*.
- [2] ANANDKUMAR, A., CHAUDHURI, K., HSU, D., KAKADE, S. M., SONG, L. and ZHANG, T. (2011). Spectral Methods for Learning Multivariate Latent Tree Structure. *Preprint, ArXiv 1107.1283*.
- [3] ANANDKUMAR, A., TAN, V. Y. F., HUANG, F. and WILLSKY, A. S. (2012a). High-Dimensional Gaussian Graphical Model Selection: Walk-Summability and Local Separation Criterion. *Accepted to J. Machine Learning Research, ArXiv 1107.1270*.
- [4] ANANDKUMAR, A., TAN, V. Y. F., HUANG, F. and WILLSKY, A. S. (2012b). High-Dimensional Structure Learning of Ising Models: Local Separation Criterion. *Accepted to Annals of Statistics*.
- [5] BAYATI, M., MONTANARI, A. and SABERI, A. (2009). Generating random graphs with large girth. In *Proc. of ACM-SIAM Symposium on Discrete Algorithms (SODA)*.
- [6] BENTO, J. and MONTANARI, A. (2009). Which Graphical Models are Difficult to Learn? In *Proc. of Neural Information Processing Systems (NIPS)*.
- [7] BHATIA, R. (2007). *Perturbation Bounds for Matrix Eigenvalues (Classics in Applied Mathematics)*. Society for Industrial and Applied Mathematics.
- [8] BLEI, D. M., NG, A. Y. and JORDAN, M. I. (2003). Latent dirichlet allocation. *J. of Machine Learning Research* **3** 993–1022.
- [9] BOGDANOV, A., MOSSEL, E. and VADHAN, S. (2008). The Complexity of Distinguishing Markov Random Fields. *Approximation, Randomization and Combinatorial Optimization. Algorithms and Techniques* 331–342.
- [10] BRÉMAUD, P. (1999). *Markov Chains: Gibbs fields, Monte Carlo simulation, and queues*. Springer.
- [11] BRESLER, G., MOSSEL, E. and SLY, A. (2008). Reconstruction of Markov Random Fields from Samples: Some Observations and Algorithms. In *Intl. workshop APPROX Approximation, Randomization and Combinatorial Optimization* 343–356. Springer.
- [12] BUNEMAN, P. (1971). The recovery of trees from measures of dissimilarity. , *Mathematics in the Archaeological and Historical Sciences (FR Hodson, DG Kendall, and P. Tautu, eds.)*.
- [13] BUNKE, H. and ALLERMANN, G. (1983). Inexact graph matching for structural pattern recognition. *Pattern Recognition Letters* **1** 245–253.
- [14] CHANDRAN, L. S. and SUBRAMANIAN, C. (2005). Girth and treewidth. *J. of combinatorial theory, Series B* **93** 23–32.
- [15] CHANDRASEKARAN, V., PARRILO, P. A. and WILLSKY, A. S. (2010a). Latent Variable Graphical Model Selection via Convex Optimization. *Preprint. Available on ArXiv*.
- [16] CHANDRASEKARAN, V., PARRILO, P. A. and WILLSKY, A. S. (2010b). Latent Variable Graphical Model Selection via Convex Optimization. *Arxiv preprint*.
- [17] CHEN, T., ZHANG, N. L. and WANG, Y. (2008). Efficient model evaluation in the search based approach to latent structure discovery. In *4th European Workshop on Probabilistic Graphical Models*.
- [18] CHOI, M. J., LIM, J. J., TORRALBA, A. and WILLSKY, A. S. (2010). Exploiting Hierarchical Context on a Large Database of Object Categories. In *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*.
- [19] CHOI, M. J., TAN, V. Y. F., ANANDKUMAR, A. and WILLSKY, A. (2011). Learning Latent Tree Graphical Models. *J. of Machine Learning Research* **12** 1771–1812.
- [20] CHUNG, F. R. K. (1997). *Spectral graph theory*. Amer Mathematical Society.
- [21] DASKALAKIS, C., MOSSEL, E. and ROCH, S. (2006). Optimal phylogenetic reconstruction. In *STOC '06: Proceedings of the thirty-eighth annual ACM symposium on Theory of computing* 159–168.

- [22] DEMBO, A. and MONTANARI, A. (2010). Ising Models on Locally Tree-like Graphs. *Annals of Applied Probability*.
- [23] DURBIN, R., EDDY, S. R., KROGH, A. and MITCHISON, G. (1999). *Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids*. Cambridge Univ. Press.
- [24] ELIDAN, G. and FRIEDMAN, N. (2005). Learning Hidden Variable Networks: The Information Bottleneck Approach. *Journal of Machine Learning Research* **6** 81-127.
- [25] ERDÖS, P. L., SZÉKELY, L. A., STEEL, M. A. and WARNO, T. J. (1999a). A few logs suffice to build (almost) all trees: Part I. *Random Structures and Algorithms* **14** 153-184.
- [26] ERDÖS, P. L., STEEL, M. A., SZÉKELY, L. A. and WARNO, T. J. (1999b). A few logs suffice to build (almost) all trees: part II. *Theoretical Computer Science* **221** 77-118.
- [27] FRANK, A. and ASUNCION, A. (2010). UCI Machine Learning Repository.
- [28] GAMBURD, A., HOORY, S., SHAHSHAHANI, M., SHALEV, A. and VIRAG, B. (2009). On the girth of random Cayley graphs. *Random Structures & Algorithms* **35** 100-117.
- [29] GEORGI, H. O. (1988). *Gibbs Measures and Phase Transitions*. Walter de Gruyter.
- [30] GRIFFITHS, R. B. (1967). Correlations in Ising ferromagnets. III. *Communications in Mathematical Physics* **6** 121-127.
- [31] JALALI, A., JOHNSON, C. and RAVIKUMAR, P. (2011). On Learning Discrete Graphical Models Using Greedy Methods. In *Proc. of NIPS*.
- [32] KARGER, D. and SREBRO, N. (2001). Learning Markov Networks: Maximum Bounded Tree-width Graphs. In *Proc. of ACM-SIAM symposium on Discrete algorithms* 392-401.
- [33] KEARNS, M. J. and VAZIRANI, U. V. (1994). *An Introduction to Computational Learning Theory*. MIT Press., Cambridge, MA.
- [34] KEMP, C. and TENENBAUM, J. B. (2008). The discovery of structural form. *Proceedings of the National Academy of Science* **105** 10687-10692.
- [35] LAURITZEN, S. L. (1996). *Graphical models*. Clarendon Press.
- [36] LAZARSFELD, P. F. and HENRY, N. W. (1968). *Latent structure analysis*. Boston: Houghton Mifflin.
- [37] MEINSHAUSEN, N. and BÜHLMANN, P. (2006). High Dimensional Graphs and Variable Selection With the Lasso. *Annals of Statistics* **34** 1436-1462.
- [38] MEZARD, M. and MONTANARI, A. (2009). *Information, physics, and computation*. Oxford University Press, USA.
- [39] MOSSEL, E. (2007). Distorted metrics on trees and phylogenetic forests. *IEEE/ACM Transactions on Computational Biology and Bioinformatics* 108-116.
- [40] MOSSEL, E. and ROCH, S. (2006). Learning nonsingular phylogenies and hidden Markov models. *The Annals of Applied Probability* **16** 583-614.
- [41] NETRAPALLI, P., BANERJEE, S., SANGHAVI, S. and SHAKKOTTAI, S. (2010). Greedy Learning of Markov Network Structure. In *Proc. of Allerton Conf. on Communication, Control and Computing*.
- [42] NEWMAN, D., BONILLA, E. V. and BUNTINE, W. (2011). Improving Topic Coherence with Regularized Topic Models. In *Proc. of NIPS*.
- [43] NEWMAN, D., KARIMI, S. and CAVEDON, L. (2009). External Evaluation of Topic Models. In *Proceedings of the 14th Australasian Computing Symposium(ACD2009)* 8.
- [44] PEARL, J. (1988). *Probabilistic Reasoning in Intelligent Systems—Networks of Plausible Inference*. Morgan Kaufmann.
- [45] RAVIKUMAR, P., WAINWRIGHT, M. J. and LAFFERTY, J. (2008). High-dimensional Ising Model Selection Using l_1 -Regularized Logistic Regression. *Annals of Statistics*.
- [46] RAVIKUMAR, P., WAINWRIGHT, M. J., RASKUTTI, G. and YU, B. (2011). High-dimensional covariance estimation by minimizing l_1 -penalized log-determinant divergence. *Electronic Journal of Statistics* **4** 935-980.
- [47] SCHWARZ, G. (1978). Estimating the dimension of a model. *Annals of Statistics* **6** 461-464.
- [48] SEMPLE, C. and STEEL, M. (2003). *Phylogenetics*. Oxford University Press.
- [49] STEEL, M. (1994). Recovering a tree from the leaf colourations it generates under a Markov model. *Applied Mathematics Letters* **7** 19-23.
- [50] WEISSMAN, T., ORDENTLICH, E., SEROUSSI, G., VERDU, S. and WEINBERGER, M. L. (2003). Inequalities for the l_1 deviation of the empirical distribution Technical Report, Hewlett-Packard Labs.
- [51] WEITZ, D. (2005). Combinatorial Criteria for Uniqueness of Gibbs Measures. *Random Structures & Algorithms* **27** 445.
- [52] ZHANG, N. L. (2004). Hierarchical Latent Class Models for Cluster Analysis. *Journal of Machine Learning Research* **5** 697-723.
- [53] ZHANG, N. L. and KOCKA, T. (2004). Efficient Learning of Hierarchical Latent Class Models. In *ICTAI*.

ELECTRICAL ENGINEERING & COMPUTER SCIENCE DEPT.,
4408 ENGINEERING HALL, IRVINE, CA, USA 92697.
E-MAIL: a.anandkumar@uci.edu; rvalluva@uci.edu